# The AI CUDA Engineer: Agentic CUDA Kernel Discovery, Optimization and Composition

Robert Tjarko Lange[1], Aaditya Prasad[1], Qi Sun[1], Maxence Faldor[1], Yujin Tang[1] and David Ha[1]

[1]Sakana AI

Recent advances in Large Language Models have driven large-scale deployment, resulting in ever-growing inference time and energy demand. While manual optimization of low-level code implementations is feasible, it is an arduous task that requires deep expertise to balance the complex interplay of algorithmic, software, and hardware bottlenecks. This report presents the first comprehensive agentic framework for *fully automatic CUDA kernel discovery and optimization*, enabling frontier large language models to perform the translation of `torch` code to CUDA kernels and then iteratively improve their runtime. We introduce THE AI CUDA ENGINEER, which acts in sequential stages. First, it translates raw PyTorch code into equivalent CUDA kernels. Next, it optimizes their runtime performance using a novel evolutionary meta-generation procedure tailored towards the CUDA ecosystem. Finally, it uses an innovation archive of discovered 'stepping stone' kernels to improve future performance on new tasks. THE AI CUDA ENGINEER can produce CUDA kernels that exceed the performance of `torch` native and compiled kernels. Out of the 250 tasks tested, THE AI CUDA ENGINEER successfully optimizes 186 tasks to a median speedup of 1.52x. For operations such as *fused 3D convolutions* or *Diagonal Matrix Multiplication*, we show runtime improvements $\geq$50x over their torch implementations. Alongside this report, we release the best discovered kernels, an accompanying dataset of all discovered kernels and an interactive webpage for exploration of the results.

## Contents

*Corresponding author(s): Robert Tjarko Lange (robert@sakana.ai)*

# 1. Introduction

The demand for computational power in machine learning has increased exponentially over the past decade, driven by the rising complexity of deep learning models and the need for large-scale data processing (Sevilla et al., 2022). Scaling laws suggest that performance improvements in language models correlate strongly with increases in compute, data, and model size (Kaplan et al., 2020). Similarly, research by Hernandez and Brown (2020) shows that the compute requirements for state-of-the-art AI models has doubled approximately every 3-4 months since 2012, vastly outpacing Moore's Law. This has led to the proliferation of specialized hardware, such as GPUs and TPUs, optimized for deep learning workloads (Jouppi et al., 2017). The emergence of foundation models, which require incredible amounts of training and inference infrastructure (Anthropic, 2023, 2024; Guo et al., 2025; Jaech et al., 2024), has exacerbated this trend and necessitates innovations in model efficiency and hardware acceleration (Bommasani et al., 2021). Additionally, the release of `DeepSeek-V3` (Liu et al., 2024) has highlighted various trade-offs between hardware and software constraints. However, the skill set required to balance such trade-offs and to engineer efficient Compute Unified Device Architecture (CUDA; Luebke, 2008) kernels is rare and highly in demand, encompassing algorithmic, hardware, and instruction set knowledge.

Simultaneously, there have been advances in automated agentic discovery (e.g. Lu et al., 2024a,b; Romera-Paredes et al., 2024) using Large Language Models (LLMs). A key advantage of LLM-driven agentic discovery is the ability of LLMs to iteratively refine hypotheses, generate experiments, and interpret results, which facilitates a closed-loop system for scientific exploration (Song et al., 2024; Yang et al., 2023). Recent work has shown LLM agents autonomously deriving new mathematical conjectures (Romera-Paredes et al., 2024), discovering novel preference optimization objectives (Lu et al., 2024a), and even writing entire Machine Learning papers (Lu et al., 2024b). This iterative refinement process mirrors an evolutionary mechanism, where LLMs act as variation-generating engines that propose new hypotheses, code implementations, and theoretical insights (Lange et al., 2024; Lehman et al., 2022; Meyerson et al., 2023). Unlike random search, however, LLMs incorporate strong algorithmic priors, allowing them to generate executable programs, testable simulations, and provable mathematical statements, significantly accelerating the search for novel solutions. Crucially, these models also integrate a notion of *human interestingness* (Faldor et al., 2024; Zhang et al., 2024b), learned through vast training corpora and human-regularized post-training. This ensures that their discoveries are not just novel but also aligned with human curiosity and problem-solving heuristics. By unifying strong coding capabilities, evolutionary search mechanisms, and alignment with human interestingness, LLMs are poised to become autonomous engines of discovery, capable of pushing the boundaries of science, engineering, and mathematics beyond the limitations of human intuition and trial-and-error experimentation.

In this paper, we set out to answer a simple question: ***Can we leverage recent agentic advances to automate the engineering and optimization of low-level inference-time computations?*** We start by investigating whether LLMs are capable of translating PyTorch code into correct and fast CUDA kernels. Next, we introduce an LLM-driven evolutionary optimization framework to directly optimize the runtime of the kernel operation. Finally, we devise a retrieval-augmented in-context prompting paradigm, which leverages an archive of previously discovered 'stepping stone' (Stanley and Lehman, 2015) kernels to aid the agent in performing translation and optimization. We compartmentalize the agentic pipeline of THE AI CUDA ENGINEER into four consecutive stages (fig. 1):
**Stage 1**: Conversion of an object-oriented PyTorch (`nn.Module`) operation into a functional version, which explicitly provides operation parameters as function inputs.
**Stage 2**: Translation of the functional PyTorch code into a numerically correct CUDA kernel. We sequentially sample individual CUDA kernel proposals and test their correctness by comparing them with the corresponding `torch` reference implementation.
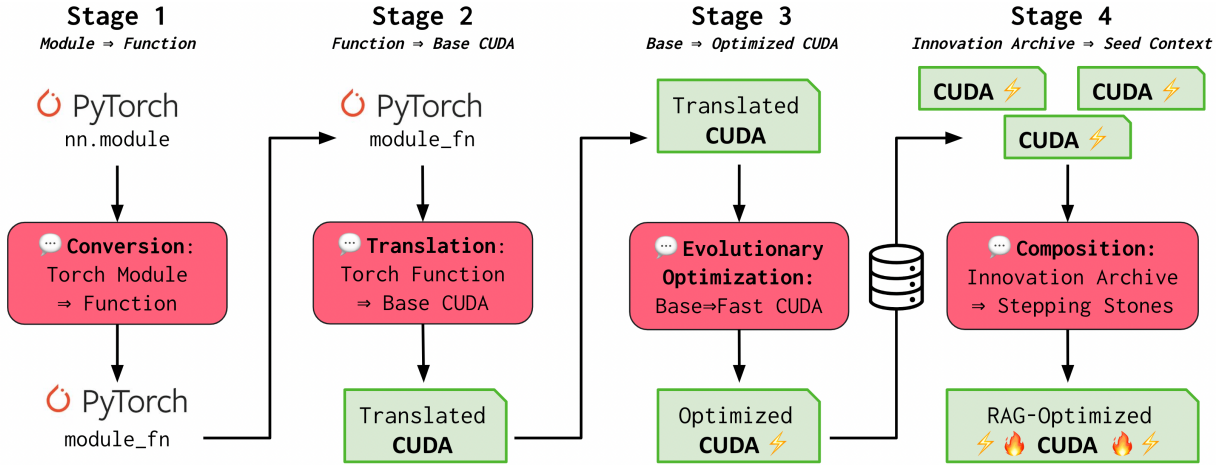
Figure 1 | **High-level overview of the AI CUDA Engineer Workflow.** In stage 1, a module definition of PyTorch code is converted into a functional version that separates out the operator parameters. In stage 2, the functional version is translated into a corresponding correct CUDA kernel, which can be loaded to replace the native PyTorch version of the operation. In stage 3, we use the translated kernel to initialize a runtime optimization process, which samples, edits, tests, and evaluates a batch of kernels in parallel. Finally, in stage 4, we leverage the discovered kernels to retrieve related ones. These are then used to enhance the in-context information for translation and/or optimization.

**Stage 3**: Evolutionary runtime optimization of the CUDA kernel while maintaining operation correctness. We sample batches of diverse CUDA kernel proposals and explicitly optimize the runtime using a combination of LLM ensembling, temperature sampling, prompting with profiling data, crossover prompting (Meyerson et al., 2023) and least-to-most sorting (Zhou et al., 2022).

**Stage 4**: Composition of in-context kernel examples given an archive of hundreds of working CUDA kernels and a target operation. We leverage code embedding-based retrieval-augmented generation (RAG, Gao et al., 2023; Lewis et al., 2020b) to seed the LLM context with useful examples.

Throughout this paper, we utilize the KernelBench (Ouyang et al., 2025) benchmark and optimize the runtime of 250 PyTorch operations. We show that THE AI CUDA ENGINEER can robustly translate PyTorch to CUDA (> 91% correctness rate), is capable of achieving significant runtime improvements (median speedup of 1.52x), and can leverage external kernel databases for improved performance. Our contributions are summarized as follows:

1. We introduce an end-to-end agentic workflow capable of translating PyTorch code to working CUDA kernels, optimizing CUDA runtime performance, and automatically fusing multiple operations (section 3). Furthermore, we construct various techniques for enhancing the consistency and performance of the pipeline including LLM ensembling, an iterative profiling feedback loop, local kernel code-editing, crossover kernel optimization, and Tensor Core support (section 8).

2. Our pipeline, THE AI CUDA ENGINEER, robustly translates more than 91% of the 250 `torch` operations and optimizes 75% of the operations better than `torch`. Furthermore, THE AI CUDA ENGINEER outperforms `torch.compile`-acceleration on 60% of the operations (section 4 and section 6).

3. We release a dataset of more than 17,000 verified kernels for operations covering a wide range of `torch` operations. These can be used for downstream applications like RAG or to fine-tune specialized models. Additionally, we release a webpage for inspecting kernels and their profiles:
   - Dataset: huggingface.co/datasets/SakanaAI/AI-CUDA-Engineer-Archive
   - Interactive WebUI: pub.sakana.ai/ai-cuda-engineer

## 2. Background

**Scaling Test-Time Compute**. Due to the increasing costs associated with training in-house models, scaling test-time compute has emerged as a strong paradigm to improve LLM outputs outside of training a new model from scratch. There are numerous options for improving samples drawn from an LLM, from the utilization of search techniques (Xie et al., 2023; Zhang et al., 2024a) to the incentivization of reasoning during post-training (Guo et al., 2025; Team et al., 2025). We focus on techniques that scale with the number of samples drawn. Parallel samples can be aggregated via voting methods (Snell et al., 2024; Trad and Chehab, 2024) or LLM debates (Du et al., 2023; Wang et al., 2025). A related strategy is evolutionary test-time compute (Berman, 2025), which aggregates samples by mutating and recombining them. Sequential samples, meanwhile, can utilize multi-turn structures (Zhou et al., 2024) to improve and correct previous responses. Parametric verifiers (Luo et al., 2024; Snell et al., 2024) can be used to score drawn samples based on their final outcome or even intermediate steps, while non-parametric (AlphaProof and AlphaGeometry teams, 2024; Li et al., 2025) verifiers range from also providing scalar scores to filtering out incorrect answers via some heuristic. Here, we focus on drawing multiple sequential streams of samples and aggregating information across them via evolutionary optimization (by selecting previous samples to place in context and recombine) guided by verifiers, such as kernel accuracy, speed, and profiling tools.

**Evolutionary Code Optimization with LLMs**. One particular flavor of test-time compute is evolutionary code optimization: the usage, mutation, and recombination of previously generated code to produce new samples. This approach has previously been used to optimize reward and preference objectives (Lu et al., 2024a; Ma et al., 2023), mathematical science code (Romera-Paredes et al., 2024), entire machine learning papers (Lu et al., 2024b), and other applications (Berman, 2025; Lange et al., 2024; Lehman et al., 2022; Meyerson et al., 2023). Through prompting, LLMs are used as recombination engines (Lange et al., 2023; Meyerson et al., 2023), and are capable of simulating crossover between diverse code snippets and the rationales that produced them. A simpler form of this technique is retrieval augmented generation (RAG, Gao et al., 2023; Lewis et al., 2020b), whereby relevant historical samples are injected into context based on embedding similarities or other filters. Here, we utilize both RAG and code-crossover to improve our kernel-optimization results.

**The CUDA Framework**. CUDA is a parallel computing platform and application programming interface (API) developed by NVIDIA for leveraging the massive parallel processing power of GPUs. It extends the C and C++ programming languages with GPU-specific extensions, allowing developers to write efficient parallelized code for tasks such as deep learning, scientific computing, and real-time graphics rendering (Nickolls et al., 2008). CUDA enables fine-grained control over GPU threads, memory hierarchies, and synchronization mechanisms, making it an essential tool for accelerating workloads that require extensive matrix and tensor computations (Kirk and Wen-Mei, 2016). Its core model consists of a hierarchy of threads organized into warps, blocks, and grids, facilitating scalable parallel execution (Garland et al., 2008). CUDA also provides libraries such as cuBLAS for linear algebra, cuDNN for deep learning, and Thrust for high-level parallel algorithms, further optimizing computational efficiency (Chetlur et al., 2014).

**KernelBench**. The KernelBench (Ouyang et al., 2025) benchmark is a set of 250 neural network tasks, defined as PyTorch modules, along with a subset of corresponding results for CUDA kernel generation across these tasks. The tasks are split into three categories denoted as levels 1, 2, and 3. Level 1 tasks are common ML primitives, such as softmax, or various matmuls. Level 2 tasks are few-step fusions of those primitives, such as network layer activations followed by a norm. Level 3 tasks represent full network architectures (e.g., a ResNet). These tasks represent a natural gradation from producing high-quality operators, to fusing them together, to assembling them into a larger neural network. In our work, we target the full suite of 250 KernelBench tasks across the three levels.

## 3. The AI CUDA Engineer: Agentic CUDA Kernel Discovery & Optimization

In this section, we provide a detailed outline of the overall workflow of THE AI CUDA ENGINEER (fig. 2). Importantly, we focus on inference and do not compute the backward pass or gradients.
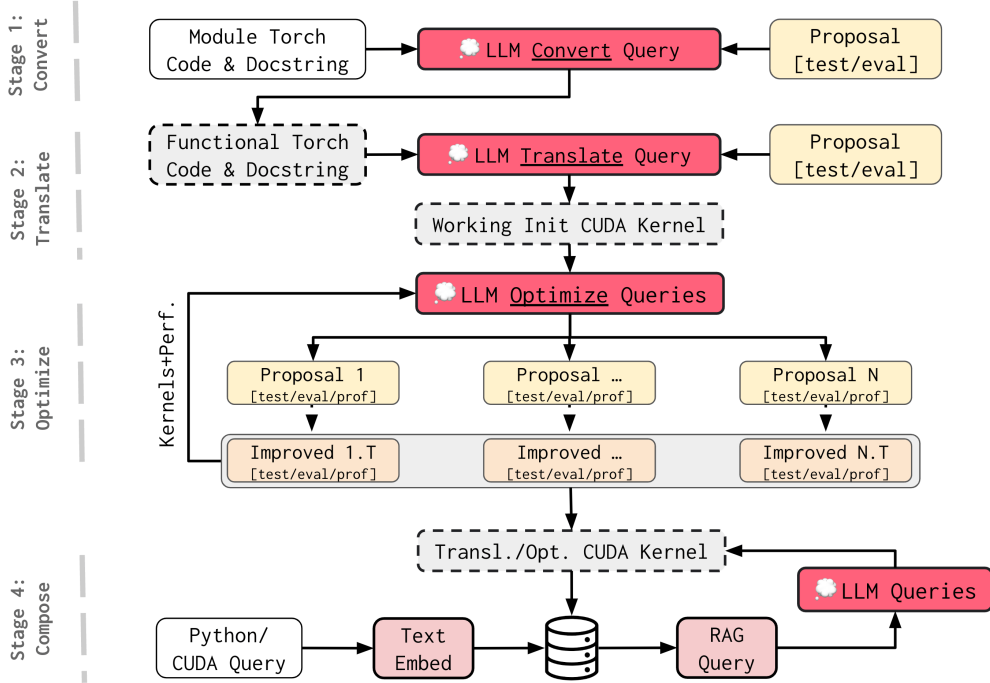


Figure 2 | **Overview of the AI CUDA Engineer.** In stage 1, user-provided `torch` code is translated into an equivalent functional representation. In stage 2, a functional PyTorch operation is translated into the corresponding CUDA kernel, which can be loaded to replace the native PyTorch version of the operation. In stage 3, we use the translated kernel to initialize a runtime optimization process, which samples, edits, tests, evaluates and profiles a batch of kernels in parallel. In stage 4, completed kernels are stored in an 'innovation archive' and used to compose in-context examples for new tasks.
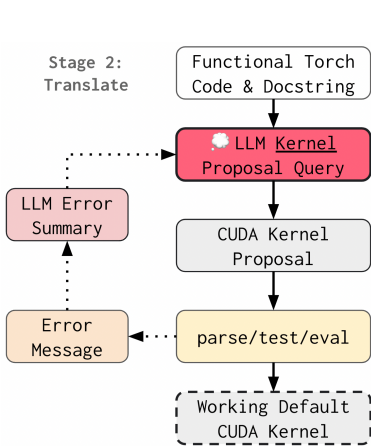


Figure 3 | **Stage 2:** Translation of functional PyTorch to CUDA

**Stage 1: Converting PyTorch Modules to Functions**.

Given a `torch` module implementation of a computation (`nn.Module`), we prompt an LLM to perform a conversion to a functional representation that explicitly separates out the function's parameters and provides a docstring description of the operation. We parse the sampled code and validate its correctness on randomly sampled inputs against the module implementation. While the approach reported in KernelBench (Ouyang et al., 2025) directly translates `torch` modules to CUDA kernels, we found this to be challenging for the LLM. The additional conversion step improves the LLMs performance when translating from `torch` to CUDA code (section 4). We hypothesize that the intermediate functional representation allows the LLM to directly reason about the underlying operation without the implicit need to infer the latent trainable parameters (e.g., weights and biases) or running summary statistics (e.g., mean and variance in normalization methods). Throughout our experiments, we observed that sequential sampling with feedback enhanced the conversion's correctness.

**Stage 2: Translating Functional PyTorch to Working CUDA**.

Next, we translate the functional PyTorch code into a working CUDA kernel. More specifically, we query an LLM with the corresponding functional `torch` implementation. After parsing the model's code output, we load the kernel using the `torch` C++ loading utilities. We check the kernel result for numerical correctness by comparing its results with the `torch` reference implementation. If the compilation fails or the computed result does not lead to a correct computation (within tolerance bounds), we summarize the error message using an additional LLM call (fig. 3). This information is fed back to the LLM before we iterate. In section 4 we show that this procedure significantly improves the robustness of correct translations as compared to best-of-N sampling (Brown et al., 2024) with the same compute budget.
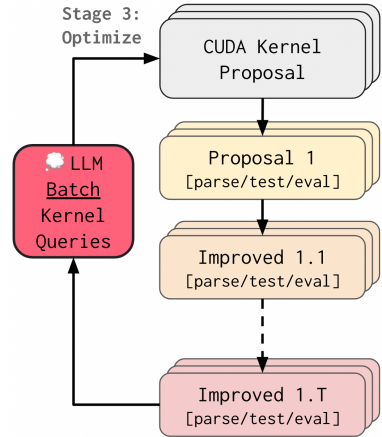


Figure 4 | **Stage 3:** Evolutionary CUDA Runtime Optimization

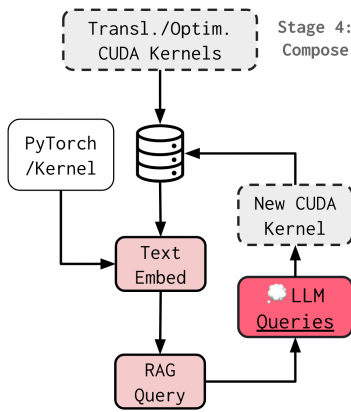**Stage 3: Evolutionary Optimization of CUDA Kernels**.



Figure 5 | **Stage 4:** RAG-Retrieval of Kernel Examples

We use the working CUDA kernel from stage 2 to kick off an evolutionary optimization process, which samples, evaluates, profiles, and edits batches of kernels to improve the runtime (fig. 4). We deploy several complementary prompting strategies:

- Given a set of previous kernel evaluations, we filter for correctness and provide a subset of up to five correct kernels sorted from slowest to fastest based on runtime speedups (Zhou et al., 2022). This least-to-most ordering enables the LLM to progressively infer optimization patterns from simpler to more sophisticated implementations.
- We sample $N = 4$ proposals from an LLM ensemble including both reasoning (o3-mini & DeepSeek-R1) and standard LLMs (Claude Sonnet 3.5 & GPT-4o). Additionally, we perform temperature sampling when applicable (Renze and Guven, 2024).
- Inspired by AlphaCode's (Li et al., 2022) prompting approach, we sample high-level recommendations (e.g., optimizing block sizes, using stride loops, etc.) to encourage diversity in model outputs. Furthermore, we introduce a novel kernel crossover prompting procedure inspired by Meyerson et al. (2023).
- For each correct kernel, we obtain the `torch`, `NCU`, and `Clang-tidy` profiling information and provide a subset of this information as additional feedback to the pipeline.
- We implement kernel code edits via function-calling. We alternate between sampling complete kernels with diverse implementation strategies and making lower-level edits.

**Stage 4: Stepping Stone Composition of CUDA Kernel Context**.
Given an 'innovation archive' of hundreds of successfully translated and optimized kernels, we leverage them as stepping stones (Stanley and Lehman, 2015). More specifically, we use retrieval-augmented generation (Lewis et al., 2020b) to obtain correct high-performing kernels from related tasks. First, we obtain text embedding vectors for all `torch` code implementations and their CUDA implementations. Given a new task, we additionally embed the new `torch` implementation and retrieve the top-$k$ closest neighboring code implementations. The corresponding CUDA kernels then supplement the in-context information used for translation or optimization. Newly discovered kernels can additionally be added to the archive in a potentially open-ended fashion.

# 4. Key Results Across AI CUDA Engineer Stages

How effectively does THE AI CUDA ENGINEER translate and optimize CUDA kernels? In this section, we present our findings across 250 CUDA kernels from KernelBench (Ouyang et al., 2025). We begin with a summary of end-to-end performance and then step through our stages one by one.

Table 1 | **End-to-End Performance Comparison Across KernelBench Levels**

| | Level 1 | | Level 2 | | Level 3 | | All Levels | |
| Baseline | Native | Compile | Native | Compile | Native | Compile | Native | Compile |
|---|---|---|---|---|---|---|---|---|
| Median Speedup (all) | 1.13 | 2.53 | 1.54 | 1.45 | 1.29 | 0.92 | 1.34 | 1.49 |
| Median Speedup (successes) | 1.52 | 3.04 | 1.63 | 1.90 | 1.42 | 1.51 | 1.52 | 2.24 |
| Successful Tasks (≥ 1x speedup) | 63 | 69 | 93 | 61 | 30 | 19 | 186 | 149 |

**End-to-End Performance:** We measure speedups as the ratio between baseline time and synthetic-kernel time on a H100, so we "break-even" at a 1x speedup and denote improvements as an optimization success. We produce 186 kernels which are as fast or faster than `torch` native implementations, and 149 kernels as fast or faster than `torch compile`[1]. *We reach a 1.34x median speedup over torch native across the 250 tasks, and a 1.52x speedup over the 186 successful tasks.*

**Functional Conversion:** We first evaluate the LLMs' ability to convert `torch` modules into parameterized function calls (stage 1). Our analysis of 250 KernelBench tasks (fig. 6) reveals distinct performance patterns across complexity levels. All tested LLMs successfully generate equivalent functional implementations for basic operations and simple fused operations (level 1, 2). However, for complex composed architectures (level 3), reasoning models (`o1-high`, `o1-preview`, `o3-mini-high`) demonstrate superior robustness, converting more than 45 tasks compared to `sonnet3.5`'s 42 tasks. Notably, `o1-high` is the only model capable of successfully translating all 50 level 3 tasks.
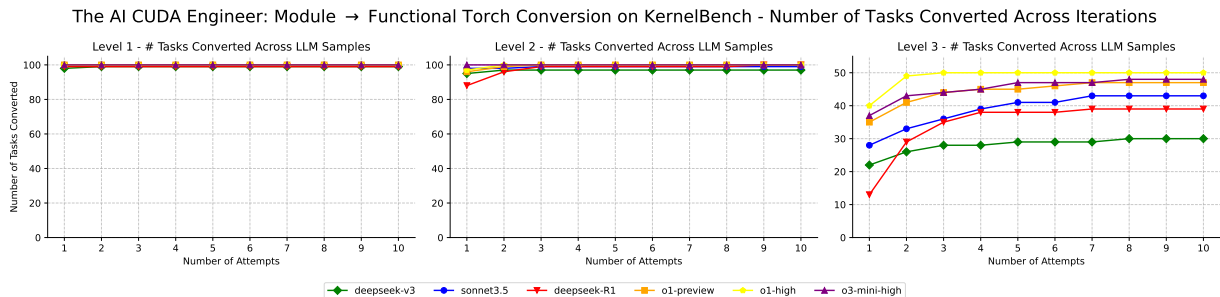


The AI CUDA Engineer: Module → Functional Torch Conversion on KernelBench - Number of Tasks Converted Across Iterations

Figure 6 | **KernelBench Conversion Results.** Most LLMs can convert the tasks into functional code at Levels 1 & 2. Notably, only `o1` completes all the tasks within a maximum of three attempts.

**CUDA Translation:** After functional conversion, we query an LLM to translate the functional `torch` implementation into an equivalent CUDA kernel (stage 2). We compare parallel proposal sampling with high temperature and sequential querying with compiler/evaluation feedback. Across all three levels, we find that incorporating compiler errors or numerical inaccuracy feedback improves the translation success rate (dashed versus solid lines in fig. 7). Increasing the number of parallel samples quickly exhibits diminishing returns (≥ 6 samples). Again, reasoning models demonstrate superior translation capabilities. They achieve a significantly higher success rate with a single sample (`pass@1`, especially for fused level 2 tasks) and are generally more capable of integrating feedback information and refining the kernel proposal successfully. Most notably, combining `o1-preview` with sequential

---

[1] `torch.compile` is not universally better than native torch, especially for the simple level 1 tasks. It is most informative to compare compile and native times/success ratios for levels 2 and 3.
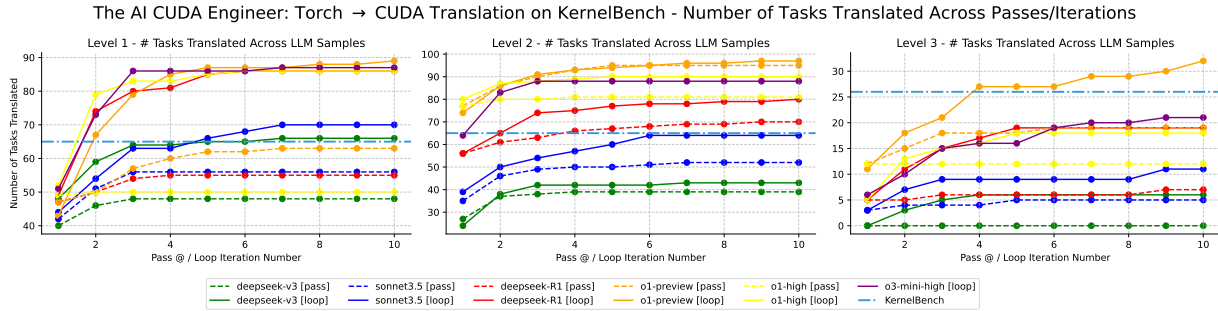
Figure 7 | **KernelBench Translation Results.** Reasoning-based LLMs (o-series and DeepSeek-R1) are capable of translating 91.6% of the considered 250 `torch` operations. Furthermore, leveraging sequential queries with error summary information significantly improves the conversion success rate.

feedback achieves a success rate of approximately 91.6%, which is significantly higher than the original results reported on the KernelBench leaderboard (approximately 62.4% as of February 2025). While the translation prompts do not explicitly target kernel runtime, we observed that these base CUDA implementations frequently outperformed their `torch` native counterparts (fig. 8, green bars).
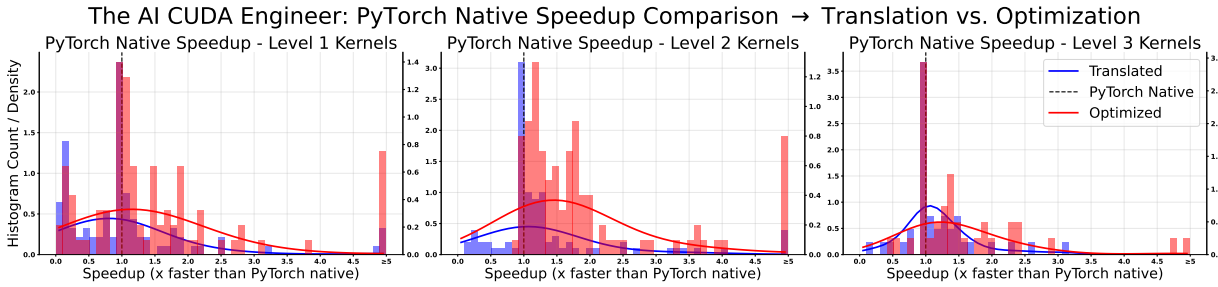


Figure 8 | **KernelBench Optimization Results.** CUDA translation alone can already result in significant runtime speedup (compared with `torch` native), evolutionary optimization unlocks additional speedups. For some tasks applying interleaving sequential proposals with code edits is effective.

**Kernel Optimization:** The translated CUDA implementation provides an initial example kernel for our downstream evolutionary optimization approach (stage 3). Using a combination of model ensembling, temperature sampling, and prompting strategies, we generate populations of kernel proposals and evaluate their correctness and runtime performance. In total, we sample 40 kernels per considered operation. While CUDA translation alone often time already achieves impressive speedups, evolutionary kernel optimization is capable of obtaining an average runtime improvement of 34% across all three KernelBench levels (fig. 8). We also observe strong speedups across a variety of diverse operations. For example, these include specialized matrix multiplications (e.g. lower triangular), Instance Normalization (Huang and Belongie, 2017), and loss functions like the Cross-Entropy Loss. We also see improvements across a plethora of operations that require substantial fusion of multiple primitives. In section 6 we provide additional details on each individual KernelBench level, and in section 8 we conduct several ablations over the design space of our optimization procedure.

**Retrieval-Augmented CUDA Kernel Translation & Optimization:** Building upon these results, we enhanced our system's capabilities through RAG. By leveraging our growing 'innovation archive' of translated and optimized kernels, RAG significantly improved both translation and optimization capabilities. We compute text embeddings of PyTorch functional code and retrieve k-nearest neighbors to seed the LLM context. This approach enabled the successful translation of 6 previously failed operations out of 15, marking a 40% improvement in translation coverage. For optimization, our analysis revealed that RAG consistently achieved the highest speedups across most tasks.

# 5. Detailed Case Study: Translating & Optimizing a Resnet18 Model

Before we present comprehensive results across all levels for THE AI CUDA ENGINEER in section 6, we highlight a specific case to illustrate our end-to-end procedure. The selected kernel `Resnet18` is part of KernelBench's level 3 and represents a famous convolution-based architecture (He et al., 2016). We present the code implementation generated by THE AI CUDA ENGINEER, discuss the detailed outputs of each stage, and highlight implementation strategies discovered during optimization.

**Conversion of `ResNet18` Module to PyTorch Functional.** As discussed in Section 3, THE AI CUDA ENGINEER first converts the PyTorch module into a "functional" representation, which separates parameters and strengthens CUDA kernel translation. Resnet18 has a complex structure with nested layers named "`BasicBlock`", making conversion non-trivial. We find that providing high-quality few-shot examples is vital for handling such compound modules. THE AI CUDA ENGINEER successfully converts Resnet18 as shown in Figure 10. For readability, we omit parts of the code with "..." and include the full code in the Appendix F. Note that the model parameters are separated into an input argument named "`params`" in the functional version. This conversion correctly preserves the model frame, transforms the "`BasicBlock`" module into "`basic_block_fn`" function, and maintains KernelBench's hyper-parameter configurations (e.g. stride, padding).

**Translation of Functional PyTorch `ResNet18` to CUDA.** Afterwards, we use the functional PyTorch implementation to seed the agentic CUDA translation. Using the proposed error feedback method (section 3), we successfully translate ResNet18 from the functional PyTorch form to a verified CUDA kernel. Listing 1 shows the resulting CUDA kernel, with code omissions similar to the previous step for readability. It is evident that THE AI CUDA ENGINEER properly follows the CUDA language grammar, preserves the correct logic, and yields the suitable PYBIND module interface at the end of the kernel implementation. THE AI CUDA ENGINEER can leverage the "torch/extension" package for facilitating the translation of complicated models like ResNet18. We would like to point out it can also generate operational code purely with basic libraries like "cuda", and "cuda_runtime" when translating basic kernels like "layernorm" in Appendix C.5. Compared to the kernel provided by KernelBench, we observed that THE AI CUDA ENGINEER is capable of translating the entire architecture to `torch` and does not only implement a subset of operations defining the full architecture in CUDA.

**Optimization of the `ResNet18` CUDA Kernel.** THE AI CUDA ENGINEER leverages 10 generations of kernels with a population size of 4 following the method introduced in Sec 3. Throughout the kernel optimization, numerous practical implementation strategies emerged, gradually achieving up to 1.44x speedup. We select three proposed kernels and present their LLM-generated thoughts in fig. 11. We also present the kernels and include the completed CUDA implementation in Appendix F. THE AI CUDA ENGINEER discovered that operator fusion (Thought 1) could significantly improve the runtime in the first few optimization steps. By writing a customized CUDA kernel that fuses the residual addition and the ReLU activation, the optimized kernel



Figure 9 | **Runtime Improvements across Optimization of ResNet18.**

achieves a 1.32x speedup gain (fig. 9). Building upon this success, THE AI CUDA ENGINEER improved the tensor data layout to better match the 2D input structure of ResNet18 (Thought 2). By leveraging shared memory for the fused operation (Thought 3) to minimize global memory accesses, our method achieved its peak performance gain of 1.44x. In addition, the monotonically increasing trend in Figure 9 suggests additional improvements could be produced by extending the search.
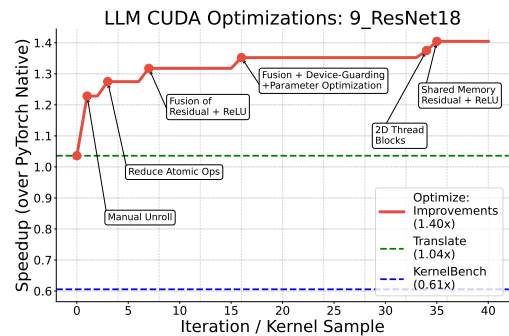
```python
...
class BasicBlock(nn.Module):
    expansion = 1

    def __init__(self, in_channels, out_channels, stride=1,
    ↪    downsample=None):
        """
        :param in_channels: Number of input channels
        :param out_channels: Number of output channels
        :param stride: Stride for the first convolutional layer
        :param downsample: Downsample layer for the shortcut
        ↪    connection
        """
        super(BasicBlock, self).__init__()
        self.conv1 = nn.Conv2d(in_channels, out_channels,
        ↪    kernel_size=3, stride=stride, padding=1, bias=False)
        self.bn1 = nn.BatchNorm2d(out_channels)
        self.relu = nn.ReLU(inplace=True)
        self.conv2 = nn.Conv2d(out_channels, out_channels,
        ↪    kernel_size=3, stride=1, padding=1, bias=False)
        self.bn2 = nn.BatchNorm2d(out_channels)
        self.downsample = downsample
        self.stride = stride

    def forward(self, x):
        """
        :param x: Input tensor, shape (batch_size, in_channels,
        ↪    height, width)
        :return: Output tensor, shape (batch_size, out_channels,
        ↪    height, width)
        """
        identity = x
        out = self.conv1(x)

        ...
        if self.downsample is not None:
            identity = self.downsample(x)
        out += identity
        out = self.relu(out)

        return out

class Model(nn.Module):
    def __init__(self, num_classes=1000):
        """
        :param num_classes: Number of output classes
        """
        super(Model, self).__init__()
        self.in_channels = 64

        self.conv1 = nn.Conv2d(3, 64, kernel_size=7, stride=2,
        ↪    padding=3, bias=False)
        self.bn1 = nn.BatchNorm2d(64)
        self.relu = nn.ReLU(inplace=True)
        self.maxpool = nn.MaxPool2d(kernel_size=3, stride=2,
        ↪    padding=1)

        self.layer1 = self._make_layer(BasicBlock, 64, 2,
        ↪    stride=1)
        ...

        self.avgpool = nn.AdaptiveAvgPool2d((1, 1))
        self.fc = nn.Linear(512 * BasicBlock.expansion,
        ↪    num_classes)

    def _make_layer(self, block, out_channels, blocks,
    ↪    stride=1):
        downsample = None
        if stride != 1 or self.in_channels != out_channels *
        ↪    block.expansion:
            downsample = nn.Sequential(
                nn.Conv2d(self.in_channels, out_channels *
                ↪    block.expansion, kernel_size=1,
                ↪    stride=stride, bias=False),
                nn.BatchNorm2d(out_channels * block.expansion),
            )

        layers = []
        layers.append(block(self.in_channels, out_channels,
        ↪    stride, downsample))
        self.in_channels = out_channels * block.expansion
        for _ in range(1, blocks):
            layers.append(block(self.in_channels, out_channels))

        return nn.Sequential(*layers)

    def forward(self, x):
        x = self.conv1(x)
        ...
        x = self.fc(x)

        return x
```

```python
...
def module_fn(x: torch.Tensor, params: nn.ParameterDict,
↪    is_training: bool):
    """
    Implements the ResNet18 module.

    Args:
        x (torch.Tensor): Input tensor, shape (batch_size,
        ↪    in_channels, height, width)
        params (nn.ParameterDict): Dictionary of parameters
        is_training (bool): Whether to use training mode

    Returns:
        torch.Tensor: Output tensor, shape (batch_size,
        ↪    num_classes)
    """
    # Initial layers
    x = F.conv2d(x, params["conv1_weight"], None, stride=2,
    ↪    padding=3)
    x = F.batch_norm(
        x,
        ...
    )
    x = F.relu(x)
    x = F.max_pool2d(x, kernel_size=3, stride=2, padding=1)

    def basic_block_fn(
        x,
        ...
    ):
        identity = x
        out = F.conv2d(x, conv1_w, conv1_b, stride=stride,
        ↪    padding=1)
        out = F.batch_norm(out, bn1_mean, bn1_var, bn1_w, bn1_b,
        ↪    is_training)
        out = F.relu(out)

        out = F.conv2d(out, conv2_w, conv2_b, stride=1,
        ↪    padding=1)
        out = F.batch_norm(out, bn2_mean, bn2_var, bn2_w, bn2_b,
        ↪    is_training)
        if downsample_conv_w is not None:
            identity = F.conv2d(x, downsample_conv_w,
            ↪    downsample_conv_b, stride=stride)
            identity = F.batch_norm(
                identity,
                ...
            )
        out += identity
        out = F.relu(out)
        return out

    # Layer blocks
    for i in range(1, 5):
        layer_name = f"layer{i}"
        for j in range(2):
            block_name = f"{layer_name}_{j}"
            stride = 2 if i > 1 and j == 0 else 1

            # Basic block parameters
            conv1_w = params[f"{block_name}_conv1_weight"]
            ...
            bn1_var = params[f"{block_name}_bn1_running_var"]

            conv2_w = params[f"{block_name}_conv2_weight"]
            ...
            bn2_var = params[f"{block_name}_bn2_running_var"]

            # Downsample parameters if they exist
            has_downsample = f"{block_name}_downsample_0_weight"
            ↪    in params
            downsample_args = {}
            if has_downsample:
                downsample_args = {
                    "downsample_conv_w": params[
                        f"{block_name}_downsample_0_weight"],
                    ...
                    "downsample_bn_var": params[
                        f"{block_name}_downsample_1_running_var"
                    ],
                }

            x = basic_block_fn(
                x,
                ...
                **downsample_args,
            )
    x = F.adaptive_avg_pool2d(x, (1, 1))
    x = torch.flatten(x, 1)
    x = F.linear(x, params["fc_weight"], params["fc_bias"])
    return x
```

Figure 10 | ResNet18 PyTorch implementations: Module-based (left) and Converted Functional (right)

```
#include <torch/extension.h>
#include <pybind11/pybind11.h>

namespace py = pybind11;

torch::Tensor basic_block_fn(
    torch::Tensor x,
    torch::Tensor conv1_w,
    ...
) {
    torch::Tensor identity = x;
    // First convolution
    x = torch::conv2d(x, conv1_w, /*bias=*/{}, /*stride=*/{stride, stride}, /*padding=*/{1, 1});
    // First batch normalization
    x = torch::batch_norm(x, /*weight=*/bn1_w, ...);
    x = torch::relu(x);
    // Second convolution
    x = torch::conv2d(x, conv2_w, /*bias=*/{}, /*stride=*/{1, 1}, /*padding=*/{1, 1});
    // Second batch normalization
    x = torch::batch_norm(x, /*weight=*/bn2_w, ...);
    // Downsample path
    if (downsample_conv_w.defined()) {
        identity = torch::conv2d(identity, downsample_conv_w, /*bias=*/{}, /*stride=*/{stride, stride});
        identity = torch::batch_norm(identity, /*weight=*/downsample_bn_w, ...);
    }
    x += identity;
    x = torch::relu(x);
    return x;
}

torch::Tensor module_fn(torch::Tensor x, py::object params_py, bool is_training) {
    auto get_param = [&](const std::string& key) -> torch::Tensor {
        return params_py.attr("__getitem__")(key.c_str()).cast<torch::Tensor>();
    };
    // Initial layers
    auto conv1_weight = get_param("conv1_weight");
    auto bn1_weight = get_param("bn1_weight");
    auto bn1_bias = get_param("bn1_bias");
    auto bn1_running_mean = get_param("bn1_running_mean");
    auto bn1_running_var = get_param("bn1_running_var");
    x = torch::conv2d(x, conv1_weight, /*bias=*/{}, /*stride=*/{2, 2}, /*padding=*/{3, 3});
    x = torch::batch_norm(
        x,
        /*weight=*/bn1_weight,
        ...
    );
    x = torch::relu(x);
    x = torch::max_pool2d(x, /*kernel_size=*/{3, 3}, /*stride=*/{2, 2}, /*padding=*/{1, 1});
    // Layer blocks
    for (int i = 1; i <= 4; ++i) {
        std::string layer_name = "layer" + std::to_string(i);
        for (int j = 0; j < 2; ++j) {
            std::string block_name = layer_name + "_" + std::to_string(j);
            int64_t stride = (i > 1 && j == 0) ? 2 : 1;
            // Basic block parameters
            auto conv1_w = get_param(block_name + "_conv1_weight");
            ...
            // Downsample parameters if they exist
            std::string downsample_conv_key = block_name + "_downsample_0_weight";
            bool has_downsample = PyMapping_HasKeyString(params_py.ptr(), downsample_conv_key.c_str()) == 1;
            torch::Tensor downsample_conv_w, downsample_bn_w, downsample_bn_b, downsample_bn_rm,
    downsample_bn_rv;
            if (has_downsample) {
                downsample_conv_w = get_param(block_name + "_downsample_0_weight");
                ...
            } else {
                downsample_conv_w = torch::Tensor();
                ...
            }
            x = basic_block_fn(
                x,
                conv1_w,
                ...
            );
        }
    }
    x = torch::adaptive_avg_pool2d(x, {1, 1});
    x = x.view({x.size(0), -1});
    auto fc_weight = get_param("fc_weight");
    auto fc_bias = get_param("fc_bias");
    x = torch::linear(x, fc_weight, fc_bias);
    return x;
}

PYBIND11_MODULE(TORCH_EXTENSION_NAME, m) {
    m.def("forward", &module_fn, "ResNet18 forward function (CUDA)");
}
```

Listing 1 | THE AI CUDA ENGINEER Translated CUDA Kernel for Resnet18.

**Thought 1: Fusing Residual addition and ReLU activation**

This version fuses the residual addition and ReLU activation into a custom CUDA kernel that uses efficient thread/block indexing. Instead of separately performing an element-wise addition and then a ReLU, the fused kernel maps each tensor element onto a thread using 1D grid indexing, thereby reducing memory accesses and kernel launch overhead.

**Thought 2: Implement 2D Thread Blocks to Match the 2D input structure**

Given the 2D nature of the convolution operations in ResNet, we can optimize the kernel by using 2D thread blocks and grid dimensions to better match the data layout. This approach reduces the overhead of index calculations and improves memory access patterns. We use a 16x16 thread block configuration which is well-suited for modern GPU architectures and the typical feature map sizes in ResNet.

**Thought 3: Shared Memory Utilization for Reduced Global Memory Access**

By utilizing shared memory for residual addition and ReLU activation, we can reduce global memory accesses. The kernel processes data in tiles using shared memory, with minimal thread synchronization points only where necessary for shared memory consistency.

```
@@ -1,8 +1,36 @@
+// Fused CUDA kernel: performs element-wise addition of x and identity, then applies ReLU
+__global__ void fused_add_relu_kernel(const float* __restrict__ x,
+                                      const float* __restrict__ identity,
+                                      float* __restrict__ out,
+                                      int numel) {
+    // Use 1D grid/block indexing for efficient mapping
+    int idx = blockIdx.x * blockDim.x + threadIdx.x;
+    if (idx < numel) {
+        float sum = x[idx] + identity[idx];
+        out[idx] = (sum > 0.0f) ? sum : 0.0f;
+    }
+}
+
+// Host wrapper to launch the fused kernel
+torch::Tensor fused_add_relu(torch::Tensor x, torch::Tensor identity) {
+    auto out = torch::empty_like(x);
+    int numel = x.numel();
+    int threads = 256;
+    int blocks = (numel + threads - 1) / threads;
+    fused_add_relu_kernel<<<blocks, threads, 0, at::cuda::getCurrentCUDAStream()>>>(
+        x.data_ptr<float>(), identity.data_ptr<float>(), out.data_ptr<float>(), numel);
+    return out;
+}
+
+// Basic block function with fused residual addition and ReLU
@@ -29,58 +57,28 @@
     x = torch::conv2d(x, conv1_w, /*bias=*/{}, /*stride=*/{stride, stride}, /*padding=*/{1, 1});
-    x += identity;
-    x = torch::relu(x);
+    // Fused residual addition and ReLU using custom CUDA kernel with proper thread/block indexing
+    x = fused_add_relu(x, identity);
     return x;
```

Figure 11 | CUDA Optimization Strategies Discovered by THE AI CUDA ENGINEER (Top) and Code Diff between initial translated kernel and implementation of Thought 1 (Bottom)

# 6. THE AI CUDA ENGINEER Evaluation on KernelBench Levels

In the following section, we take an in-depth look at the results obtained by THE AI CUDA ENGINEER on the three individual KernelBench (Ouyang et al., 2025) sweep levels.

## 6.1. KernelBench Level 1 Results: Optimizing Basic Operations

Level 1 of KernelBench consists of 100 PyTorch tasks that can roughly be grouped as follows: 18 tasks covering various matrix multiplication algorithms, 13 tasks for elementwise activation functions, 17 tasks for normalization methods, 5 tasks for pooling techniques, 6 tasks for all-reduce operations, 33 tasks for convolution methods, and 12 tasks covering common loss functions. Table 2 summarizes the translation success of various frontier LLMs across two main settings: First, we evaluate simple best-of-N parallel sampling of translation CUDA kernel proposals (Ehrlich et al., 2025). Second, we use a simple sequential feedback loop with error message information (as outlined in section 3). We find that for both settings `o1-preview` with 16384 completion tokens outperforms all other considered LLMs. Furthermore, incorporating error verification leads to a substantial improvement in translation success. While `o1-preview` without feedback is only capable of translating 63 out of 100 tasks, the same model with feedback obtains 89 successful translations. This represents a stronger process than previously reported on the KernelBench leaderboard (65 out of 100 tasks).

Table 2 | **Evaluation of AI CUDA Engineer PyTorch Translation for Level 1 KernelBench Tasks.**

| Level 1 [Basic Ops] | # Translated Operations | Improved over torch | Improved over compile | Improved over KernelBench | Torch P50 Speedup | Torch P75 Speedup | Torch Max Speedup |
|---|---|---|---|---|---|---|---|
| deepseek-v3 [pass@10] | 48 | 8 | 8 | 18 | 0.27 | 0.93 | 23.93 |
| sonnet3.5 [pass@10] | 56 | 13 | 13 | 20 | 0.63 | 0.97 | 29.82 |
| deepseek-R1 [pass@10] | 55 | 14 | 13 | 21 | 0.42 | 1.00 | 23.50 |
| o1-preview [pass@10] | 63 | 13 | 11 | 21 | 0.45 | 0.96 | 34.68 |
| o1-high [pass@10] | 50 | 11 | 10 | 21 | 0.37 | 0.97 | 39.63 |
| deepseek-v3 [loop@10] | 66 | 13 | 12 | 23 | 0.63 | 0.95 | 24.13 |
| sonnet3.5 [loop@10] | 70 | 16 | 16 | 23 | 0.70 | 0.98 | 29.51 |
| deepseek-R1 [loop@10] | 86 | 17 | 17 | 22 | 0.27 | 0.91 | 34.04 |
| o1-preview [loop@10] | 89 | 13 | 12 | 25 | 0.43 | 0.93 | 35.58 |
| o1-high [loop@10] | 86 | 17 | 15 | 23 | 0.31 | 0.94 | 29.83 |
| o3-mini-high [loop@10] | 87 | 17 | 15 | 23 | 0.45 | 0.96 | 29.54 |
| Union - Max All | 91 | 38 | 36 | 46 | 0.97 | 1.16 | 39.63 |

Interestingly, while we don't explicitly prompt the LLM to provide a runtime improvement, we observe that translating `torch` native code into equivalent CUDA code often results in strong speedups. When inspecting the best kernels across all translation settings, we find that 25% of successful translations led to an improvement of more than 16% (versus `torch` native). Furthermore, 36 kernels even outperform `torch`-compiled versions of the operation.
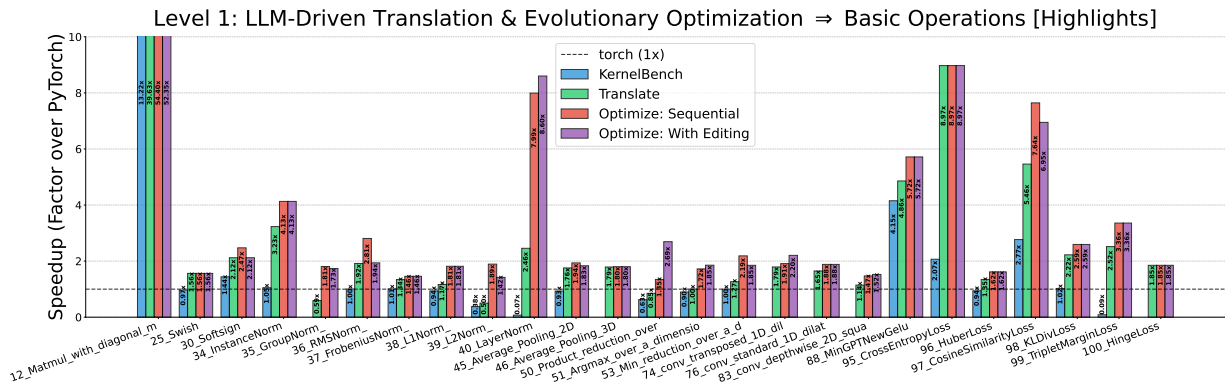


Figure 12 | **Level 1 Translation and Evolutionary Optimization Runtime Improvement Highlights.**

We can further improve performance by using THE AI CUDA ENGINEER to directly optimize the

runtime. Figure 12 and table 3 highlights a subset of improved kernels from level 1. In general, we observed that it was challenging to improve the matrix multiplication kernels. Such operations often rely on highly optimized and proprietary CUTLAS implementations, which aren't represented in the LLM training data.

Table 3 | **Highlighted AI CUDA Engineer Kernels in Level 1 of KernelBench.**

| Task | Problem Name | Translate Speedup vs. Native | Translate Speedup vs. Compile | Optimize Speedup vs. Native | Optimize Speedup vs. Compile |
|---|---|---|---|---|---|
| 12 | 12_Matmul_with_diagonal_matrices_ | 39.63 | 40.40 | 56.76 | 58.28 |
| 95 | 95_CrossEntropyLoss | 8.97 | 2.45 | 5.16 | 5.17 |
| 40 | 40_LayerNorm | 2.46 | 0.20 | 7.84 | 0.69 |
| 97 | 97_CosineSimilarityLoss | 5.46 | 3.91 | 7.0 | 5.51 |

## 6.2. KernelBench Level 2 Results: Optimizing Fused Operations

Level 2 of the KernelBench task sweep contains 100 composed operation blocks, which require the fusion of multiple primitives. These can be grouped as follows: 18 tasks consider 2D convolutions with various activation functions and simple arithmetic operations; 14 tasks implement 2D convolutions with normalization or pooling; 24 tasks combine 3D convolutions with activation and simple operations; 7 tasks synthesize 3D convolutions with normalization or pooling; 22 for generalized matrix multiplication with various residual operations activations, or arithmetic operations; and 15 tasks consider generalized matrix multiplication with normalization or pooling. Table 4 details translation rates across Level 2 tasks while Table 5 showcases optimization results.

Table 4 | **Evaluation of AI CUDA Engineer PyTorch Translation for Level 2 KernelBench Tasks.**

| Level 2 [Fused Ops] | # Translated Operations | Improved over torch | Improved over compile | Improved over KernelBench | Torch P50 Speedup | Torch P75 Speedup | Torch Max Speedup |
|---|---|---|---|---|---|---|---|
| deepseek-v3 [pass@10] | 39 | 14 | 7 | 9 | 0.49 | 1.08 | 8.14 |
| sonnet3.5 [pass@10] | 52 | 20 | 14 | 13 | 0.39 | 1.27 | 10.81 |
| deepseek-R1 [pass@10] | 70 | 23 | 16 | 17 | 0.37 | 1.06 | 59.19 |
| o1-preview [pass@10] | 95 | 42 | 13 | 22 | 1.00 | 1.01 | 3.35 |
| o1-high [pass@10] | 81 | 29 | 17 | 19 | 0.87 | 1.00 | 4.34 |
| deepseek-v3 [loop@10] | 43 | 15 | 10 | 12 | 0.96 | 1.11 | 3.35 |
| sonnet3.5 [loop@10] | 64 | 17 | 15 | 11 | 0.31 | 1.03 | 10.16 |
| deepseek-R1 [loop@10] | 80 | 21 | 18 | 14 | 0.26 | 1.00 | 112.45 |
| o1-preview [loop@10] | 97 | 45 | 19 | 26 | 1.00 | 1.04 | 3.38 |
| o1-high [loop@10] | 90 | 24 | 13 | 15 | 0.51 | 1.00 | 6.78 |
| o3-mini-high [loop@10] | 88 | 20 | 20 | 17 | 0.23 | 0.93 | 8.64 |
| Union - Max All | 98 | 72 | 40 | 41 | 1.05 | 1.31 | 112.45 |

Again, we observe that for translation `o1-preview` achieves the highest success rate. We hypothesize, that the performance difference between `o1-high` and `o1-preview` may result from potential post-training and quantization. Furthermore, we obtained 98 successful CUDA translations as compared to 65 for the KernelBench approach. Translation alone, again, achieves a speedup improvement over 40 compiled `torch` implementations. After optimization, that number jumps to 61, with a median speedup of 1.9x over those kernels. The highlighted kernels in table 5 showcase that THE AI CUDA ENGINEER is capable of obtaining improvements for a broad range of computations including various composed matrix multiplication kernels, 2D as well as 3D convolutions. These individual and aggregated results highlight that THE AI CUDA ENGINEER is capable of fusing multiple operations effectively. Finally, we observe that interleaving batch proposal sampling with code editing improves the runtime obtained by THE AI CUDA ENGINEER for some tasks but not uniformly. Hence, the optimal trade-off between sampling diverse ideas and locally improving implementation details appears to be task-sensitive.
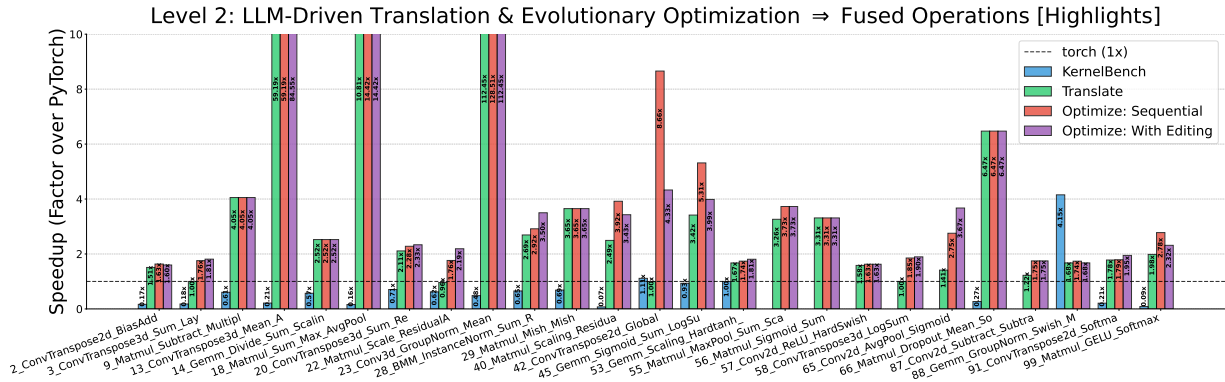
Figure 13 | **Level 2 Translation and Evolutionary Optimization Runtime Improvement Highlights.**
Table 5 | **Highlighted AI CUDA Engineer Kernels in Level 2 of KernelBench.**

| Task | Problem Name | Translate Speedup vs. Native | Translate Speedup vs. Compile | Optimize Speedup vs. Native | Optimize Speedup vs. Compile |
|------|--------------|------------------------------|-------------------------------|------------------------------|-------------------------------|
| 23 | 23_Conv3d_GroupNorm_Mean | 112.45 | 72.04 | 84.67 | 73.59 |
| 13 | 13_ConvTranspose3d_Mean_Add_Softmax_Tanh | 59.19 | 46.77 | 84.01 | 65.83 |
| 18 | 18_Matmul_Sum_Max_AvgPool_LogSumExp_LogS | 10.81 | 1.98 | 21.08 | 13.99 |
| 66 | 66_Matmul_Dropout_Mean_Softmax | 6.47 | 11.33 | 7.21 | 13.47 |

## 6.3. KernelBench Level 3 Results: Optimizing Complex Architectures

Finally, level 3 of the KernelBench task sweep encompasses a set of 50 popular neural network architecture blocks, which often entail hundreds of lines of `torch` code. As seen in table 6, our pipeline achieves successful translation over forty of the fifty level 3 tasks. Interestingly, we find the `o1-preview` with error feedback scales well for increased translation budgets (fig. 7). `DeepSeek-v3`, on the other hand, is not able to translate a single level 3 kernel. After optimization, we have kernels as fast or faster than torch compile on nineteen of the fifty tasks, with a median speedup of 1.51x. We highlight some of these kernels in table 7.



Figure 14 | **Level 3 Translation and Evolutionary Optimization Runtime Improvement Highlights.**

These highlighted kernels include both recurrent neural network implementations and classic vision architectures. Generally, speaking the observed maximal speedups obtained for level 3 tend to be smaller (e.g. 3x) as compared to levels 1 and 2. This may likely be due to increased manual optimization for popular architectures and the increased context and output length required for such involved architectures. Finally, we again observed that the previously reported KernelBench (Ouyang et al., 2025) is oftentimes only capable of factoring out small subsets of the overall architectures into working CUDA code. THE AI CUDA ENGINEER, on the other hand, is capable of obtaining CUDA

Table 6 | **Evaluation of AI CUDA Engineer PyTorch Translation for Level 3 KernelBench Tasks.**

| Level 3 [Complex Models] | # Translated Operations | Improved over torch | Improved over compile | Improved over KernelBench | Torch P50 Speedup | Torch P75 Speedup | Torch Max Speedup |
|---|---|---|---|---|---|---|---|
| deepseek-v3 [pass@10] | 0 | 0 | 0 | 0 | nan | nan | nan |
| sonnet3.5 [pass@10] | 5 | 0 | 0 | 0 | 0.10 | 0.23 | 0.75 |
| deepseek-R1 [pass@10] | 7 | 0 | 0 | 1 | 0.38 | 0.56 | 0.91 |
| o1-preview [pass@10] | 19 | 4 | 0 | 5 | 0.99 | 1.00 | 1.70 |
| o1-high [pass@10] | 12 | 4 | 1 | 4 | 0.93 | 1.00 | 1.57 |
| deepseek-v3 [loop@10] | 6 | 2 | 2 | 1 | 0.22 | 1.75 | 3.03 |
| sonnet3.5 [loop@10] | 11 | 3 | 1 | 3 | 1.00 | 1.00 | 3.01 |
| deepseek-R1 [loop@10] | 19 | 11 | 6 | 10 | 1.00 | 1.46 | 3.10 |
| o1-preview [loop@10] | 32 | 13 | 4 | 14 | 1.00 | 1.05 | 1.62 |
| o1-high [loop@10] | 18 | 8 | 0 | 3 | 0.55 | 1.08 | 2.58 |
| o3-mini-high [loop@10] | 21 | 2 | 2 | 3 | 0.11 | 0.38 | 1.90 |
| Union - Max All | 40 | 26 | 9 | 18 | 1.01 | 1.42 | 3.10 |

code that implements the entire computation graph.

Table 7 | **Highlighted AI CUDA Engineer Kernels in Level 3 of KernelBench.**

| Task | Problem Name | Translate Speedup vs. Native | Translate Speedup vs. Compile | Optimize Speedup vs. Native | Optimize Speedup vs. Compile |
|---|---|---|---|---|---|
| 34 | 34_VanillaRNNHidden | 1.01 | 2.31 | 7.02 | 11.61 |
| 24 | 24_EfficientNetB2 | 1.49 | 1.02 | 2.37 | 1.63 |
| 4 | 4_LeNet5 | 2.25 | 1.38 | 2.38 | 1.47 |

# 7. Retrieval-Augmented CUDA Kernel Translation & Optimization

**Translation** Given a dataset of kernels discovered by THE AI CUDA ENGINEER, we leverage RAG to improve the translation capabilities of the system, leading to self-improvement. For each kernel in our archive, we first compute text embeddings of their PyTorch functional code (algorithm 1). When translating a new operation, we compute its embedding and retrieve the $k$

Table 8 | **Six CUDA kernel translations enabled with RAG;**

| Operation | Level | Speedup |
|---|---|---|
| Conv3D (Standard, Square) | 1 | 1.0 |
| Conv3D (Standard, Asymmetric) | 1 | 0.02 |
| Conv3D (Transposed, Square) | 1 | 1.0 |
| Conv3D (Transposed, Asymmetric) | 1 | 0.03 |
| Conv3D (Transposed, Padded) | 1 | 0.01 |
| Conv3D-GroupNorm-Min-Clamp-Dropout | 2 | 1.0 |

most similar PyTorch implementations along with their corresponding CUDA kernels. These examples are then provided as in-context demonstrations to guide the translation process. Out of 15 previously failed tasks, RAG enabled the successful translation of 6 operations, a 40% improvement. As shown in Table 8, while three operations maintained performance parity with their PyTorch counterparts, the remaining three exhibited lower performance. However, even for these slower kernels, RAG successfully provided numerically accurate implementations that can serve as starting points for the optimization pipeline described in Section 3. Notably, the successful translation of a Level 2 fusion task demonstrates RAG's potential for handling more complex operations.

**Optimization** We also investigate the impact of RAG on CUDA kernel optimization (stage 3). Similar to the translation setup in Algorithm 1, we compute text embeddings for each kernel in our archive, but now we include optimized kernel codes for similar functions. We compare four different settings: optimization without RAG, RAG using a dataset of optimized kernels, RAG using a dataset of translated (but not optimized) kernels, and RAG with randomly selected examples instead of k-nearest neighbors.

---

**Algorithm 1** Retrieval-Augmented CUDA Kernel Translation & Optimization

---

**Require:** An archive of previously translated CUDA kernels $\mathcal{A} = \{s_1, \ldots, s_{|\mathcal{D}|}\}$, an embedder, a
   number of neighbors $k$.
 1: **for all** kernels $s \in \mathcal{A}$ **do**
 2:     embed[$s$] ← Text embedding of PyTorch functional code
 3: **end for**
 4: embed ← Text embedding of current PyTorch functional code to translate
 5: Get $k$-nearest neighbors of current PyTorch functional or CUDA kernel code in embedding space.
 6: Add $k$ neighbors (PyTorch functional, CUDA kernel code) pairs to the prompt and query LLM.

---

Figure 15 presents the runtime improvements across these settings for a diverse set of operations. Our analysis reveals that RAG-based on translated kernels consistently achieves the highest speedups across most tasks, with the notable exception of task 42 (Max Pooling 2D) where the baseline without RAG performs best. The effectiveness of nearest-neighbor selection is validated by the consistently poor performance of random examples, which never achieves the best performance for any task.



Figure 15 | **CUDA Kernel Optimization Analysis for Different RAG Settings.**

Using a dataset of optimized CUDA kernels for RAG does not improve performance compared to using translated but unoptimized kernels. This suggests that leveraging a diverse set of basic implementation strategies can be more beneficial than using highly specialized optimizations. Simpler translated kernels may provide clearer patterns for the LLM to build upon during optimization, while already-optimized kernels contain complex techniques that are harder to adapt to new contexts.

# 8. Ablation Studies of THE AI CUDA ENGINEER Components

In this section, we conduct ablation studies to understand the impact of different components in our optimization pipeline. First, we investigate how model diversity affects the quality of optimized kernels by varying the number of LLMs used in the optimization process (from 1 to 4 models). Second, we investigate the usefulness of using kernel profiling data as part of the prompting strategy. Finally, we examine the effectiveness of our crossover operation strategy in generating improved kernel implementations. Our crossover approach draws inspiration from genetic algorithms but is specifically tailored for CUDA kernel optimization, and involves the selection of high-performing parent kernels which are fused by LLMs to maintain correctness while merging features from both parents.
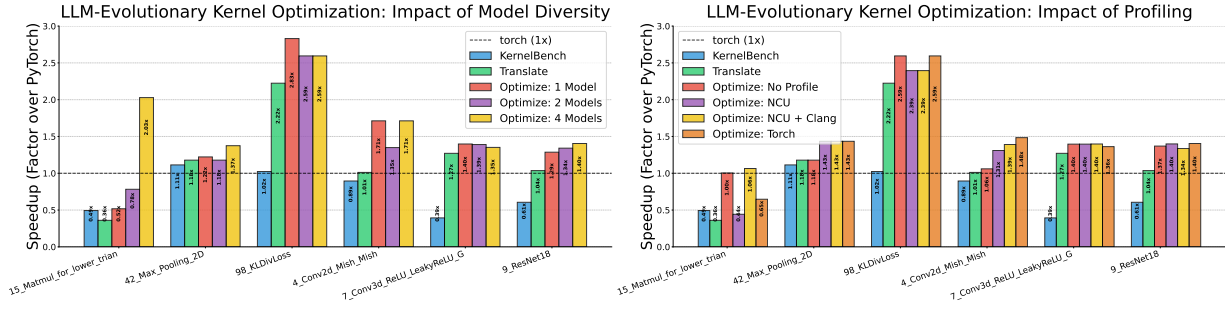


Figure 16 | **Level 3 Optimization Ablations for LLM Model Diversity & Profiling Data.**

**Impact of Model Diversity** To evaluate the effect of model diversity on kernel optimization performance, we compare three distinct configurations:

1. **Single LLM**: Using only one model (`claude-3.5-sonnet`)
2. **Two LLMs**: Combining `claude-3.5-sonnet` with `textttgpt-4o`
3. **Four LLMs**: Full ensemble using additionally `o3-mini`, and `deepseek-R1`

On average, increasing model diversity leads to improved optimization outcomes across different tasks (fig. 16). The four-model ensemble configuration achieves superior performance in most cases, with particularly notable improvements in specific tasks. For instance, in the 15_Matmul_for_lower_trian task, the four-model ensemble achieves a remarkable 2.0x speedup, significantly outperforming both the single-model (0.75x) and two-model (0.85x) configurations. For the 23_Conv3d_GroupNorm_Mean task, the four-model ensemble maintains consistent performance improvements throughout optimization. We hypothesize that the improved performance with multiple models stems from several factors:

- Diverse implementation strategies arising from different model training distributions
- Complementary optimization techniques learned by different model architectures
- Increased exploration of the solution space through varied model perspectives

**Prompting with Profiling Data** Next, we investigated the impact of including profiling data in the kernel batch sampling prompt of THE AI CUDA ENGINEER. We hypothesize that this feedback allows the LLM to reason about implementation details such as sub-operation-specific memory consumption and runtime. Again, we investigate 4 separate settings: Using no profiling information, using only NVIDIA Compute Utility (NCU) profiling data, using both NCU as well as `Clang-Tidy` information, and using only `torch` profiling data. NCU profiles provide meta-data about execution timing, resource utilization, memory metrics as well as instruction analysis. `Clang-Tidy`, on the other hand, yields a static code analysis of the source code without explicit kernel compilation or evaluation. We format the obtained raw profile before extending the context of the LLM. Including profiling data, in general, helps THE AI CUDA ENGINEER to obtain performance improvements (fig. 16). Different tasks benefit from different profiles. E.g. the max pooling kernel discovery benefits from all individual profiling settings. The fused 2D convolution task, on the other hand, benefits most from the `torch` profile. On average `torch` profiling data yields robust speedup gains.

**Kernel Crossover Prompting**    We develop a novel evolutionary approach that leverages learned code embeddings to identify and combine complementary optimization strategies from different kernels (outlined in algorithm 2). Thereby, THE AI CUDA ENGINEER is capable of 'crossing over' different CUDA kernels and combining complementary strengths. We first compute text embeddings for all kernels in our archive and project them into a lower-dimensional space. These embeddings capture semantic similarities between different implementation strategies. We then cluster the kernels and select the best-performing parents from each cluster for recombination. To evaluate the effect of

---

**Algorithm 2** Language Model CUDA Kernel Crossover

---

**Require:** An archive of previously evaluated CUDA kernels $\mathcal{A} = \{s_1, \ldots, s_{|\mathcal{D}|}\}$.
**Require:** Number of latent dimensions $Z = 2$ for code embeddings
**Require:** Number of clusters $K = 3$ to group kernels.
  1: **for all** kernels $s \in \mathcal{A}$ **do**
  2:     embed$[s] \leftarrow$ Text embedding of kernel code script
  3: **end for**
  4: Compute low-dimensional projection of all text embeddings into $Z$-dimensional space.
  5: Cluster projected kernel scripts into $K$ clusters
  6: Select the best performing (runtime) kernel as parents from each cluster
  7: Construct crossover prompts incentivizing the combination of kernel optimization strategies

---

crossover on kernel optimization performance, we compare three distinct configurations:

1. **No crossover**: Standard optimization without kernel recombination
2. **Frequent crossover**: Performing crossover every 2 generations.
3. **Moderate crossover**: Performing crossover every 4 generations.



Figure 17 | **Level 3 Optimization Ablation for Kernel Crossover Prompting.**

The variant without crossover consistently performs worst across all evaluated tasks (fig. 17). Both crossover configurations (every 2 or 4 batches) generally achieve superior performance, although there is no clear winner between these two settings. Most notably, for the `15_Matmul_for_lower_trian` task, the frequent crossover variant (every 2 batches) achieves a remarkable speedup of 4.51x, substantially outperforming the second-best configuration which only reaches 0.65x speedup.

## 9. The AI CUDA Engineer Archive: A Verified CUDA Kernel Dataset

Along with this paper, we release The AI CUDA Engineer archive, a dataset consisting of approximately 30,000 CUDA kernels generated by The AI CUDA Engineer. It is released under the CC-By-4.0 license and can be accessed via HuggingFace[2] and interactively visualized here[3]. The dataset includes a torch reference implementation, torch, NCU and Clang-tidy profiling data, multiple kernels per task, error messages and speedup scores against torch native and compile runtimes. We envision that this dataset can enable post-training of open-source models to perform better CUDA-enabling modules. Table 9 provides key summary statistics of the dataset:

Table 9 | **Summary Statistics of The AI CUDA Engineer Archive.**

| Level | Total Kernels | Correct Kernels | In-Correct Kernels | # Kernels faster Torch Native | # Kernels faster Torch Compile | # Error (Eval or Compilation) | # Error (Accuracy) |
|---|---|---|---|---|---|---|---|
| 1 | 12,157 | 7,222 | 4,935 | 2,779 | 4,339 | 3,336 | 1,599 |
| 2 | 12,938 | 6,988 | 5,950 | 4,750 | 3,235 | 4,076 | 1,874 |
| 3 | 5,520 | 3,271 | 2,249 | 1,161 | 521 | 2,660 | 611 |
| Total | 30,615 | 17,481 | 13,134 | 8,690 | 8,086 | 10,072 | 4084 |

Inspired by Lu et al., fig. 18 visualizes the different generated kernels. We embed all kernels using OpenAI's text-embedding-3-small and compute 2-dimensional projections using tSNE (Van der Maaten and Hinton, 2008), which we then cluster using DBSCAN (Ester et al., 1996). Afterwards, we pick 50 kernels for each cluster and query o3-mini to construct cluster summary descriptions.



Figure 18 | **The AI CUDA Engineer Archive Kernel Summary Visualization.** tSNE projections and DBSCAN clustering of CUDA kernel code contained in the archive. The dataset consists of more than 15,000 verified kernels which cluster into various task groups and kernel optimization strategies.

The clusters group tasks (e.g., matrix multiplication, convolution, etc.) conditioned on various kernel implementation strategies (e.g., fusion, shared memory, vectorization, etc.). This can potentially allow for targeted fine-tuning of base language models. Finally, the error messages can be used to train reasoning models without the need for online compilation and kernel evaluation.

---

[2] https://huggingface.co/datasets/SakanaAI/AI-CUDA-Engineer-Archive
[3] https://pub.sakana.ai/ai-cuda-engineer

# 10. Related Work

**Modern GPU Programming Frameworks and Automated Kernel Optimization**. Recent advances in GPU programming frameworks have fostered an ecosystem that emphasizes both high performance and developer productivity. Frameworks like Triton (Tillet et al., 2019), developed by OpenAI, offer a Python-based interface that allows developers to write efficient, low-level GPU kernels without delving into the intricate details of CUDA or OpenCL. Emerging tools such as ThunderKittens (Spector et al., 2024) focus on enhanced human usability and seamless integration with popular machine learning libraries, further simplifying complex workflows. Additionally, there are innovative methods for automating the writing of CUDA kernels. Domain-specific languages and auto-tuning frameworks—such as TVM (Chen et al., 2018) and its auto-scheduler Ansor (Zheng et al., 2020) - along with template metaprogramming in C++ and Python-based DSLs, enable the automatic generation of highly optimized kernels tailored to specific hardware and workload requirements. METR (2025) leverage LLMs for CUDA kernel generation on a modified KernelBench sweep. They report average speedups, which make it hard to assess the reported robustness. Finally, Chen et al. (2025) report promising results for optimizing attention kernels using DeepSeek-R1. Unlike THE AI CUDA ENGINEER these approaches do not detail an end-to-end system which utilizes the evolutionary kernel optimization paradigm, crossover-prompting, model ensembling or profiling data. Furthermore, we release the entire kernel dataset discovered by THE AI CUDA ENGINEER.

**Scientific Discovery with LLMs**. As LLMs become more capable, interest has grown in using them to aid scientific discovery. Prior works have shown LLMs' effectiveness in paper review (Liu and Shah, 2023; Zhuang et al., 2025), idea generation (Dasgupta et al., 2024; Si et al., 2024; Susnjak et al., 2025), and research synthesis (Agarwal et al., 2024; Ali et al., 2024). Other works have gone further in integrating LLMs across the scientific process, having them propose, implement, and evaluate preference optimization objectives (Lu et al., 2024a) as well as entire scientific experiments (Lu et al., 2024b). The latter work, The AI Scientist, is additionally capable of visualizing its results, writing a paper, and preparing it for review.

**Software Engineering with LLMs**. Another domain LLMs have shown promising results in is software engineering. Common benchmarks test standalone programming questions (Austin et al., 2021; METR, 2025; Quan et al., 2025) as well as traditional software engineering tasks, such as completing GitHub issues (Jimenez et al., 2024). A growing subset of work in this field has also begun to target kernel writing (Ouyang et al., 2025; Zhang et al., 2025) in Architecture Specific Programming Languages such as CUDA, since efficiency improvements from these kernels can be incredibly valuable and skilled human kernel engineers are in high demand. Different methods diverge in how they tackle these applications. Some works train LLMs to become better coders (Hui et al., 2024; Roziere et al., 2023; Zhu et al., 2024), while others leverage test-time compute strategies to improve model outputs. There is also variation in the output form factor, which ranges from full files to diffs that can be iteratively applied, as in Gauthier (2024).

**Test-Time Compute Scaling**. One dimension of test-time scaling is increasing the effort per sample, whether by training (Guo et al., 2025; Team et al., 2025), prompting longer model "thought" (Muennighoff et al., 2025; Wei et al., 2022), scaling search across trajectories (Snell et al., 2024; Xie et al., 2023; Zhang et al., 2024a), retrieving context (Lewis et al., 2020a), or using multi-turn responses (Zhou et al., 2024). Alternatively, more samples can be generated and aggregated via voting (Snell et al., 2024; Trad and Chehab, 2024) or LLM communication (Du et al., 2023; Wang et al., 2025). Evolutionary test-time compute (Berman, 2025; Lu et al., 2024a) mutates and recombines samples, while generating parallel outputs with different models (Du et al., 2023) or prompts (AlphaProof and AlphaGeometry teams, 2024) encourages diversity.

# 11. Limitations & Ethical Considerations

**Focus on inference runtime & forward pass computation**  Our work focuses exclusively on optimizing inference-time computations and forward pass operations. While this scope allows for meaningful runtime improvements, it represents only a subset of the complete computational deep learning pipeline. The generation of coherent forward and backward passes introduces significant additional complexity, particularly regarding the storage and management of activations required for computing Jacobian-vector products. Future work could explore extending THE AI CUDA ENGINEER to handle these more complex scenarios.

**Benchmark Quality**  We use the KernelBench tasks as-is to maintain a simple point of comparison against KernelBench's (Ouyang et al., 2025) results. A recent related work (METR, 2025) filtered down the set of tasks to remove trivial or abnormal examples. We do not do this, but we sanity-check our top kernels to ensure accuracy over a range of inputs and the qualitative absence of exploitative/trivial CUDA code. Still, more work has to be done to improve benchmark and response quality, and tasks and kernels slated for deployment must be rigorously checked.

**Challenges & Common Failure Modes**  Several technical challenges and common failure modes emerged during our investigation:

- Combining evolutionary optimization with LLMs is powerful but can also find ways to trick the verification sandbox. We are fortunate that Twitter user `@main_horse` helped test our CUDA kernels, to identify that THE AI CUDA ENGINEER had found a way to 'cheat'. The system had found a memory exploit in the evaluation code which, in a small percentage of cases, allowed it to avoid checking for correctness. We have consequently made the evaluation harness more robust to eliminate such loopholes.
- We observed limitations in frontier LLMs' ability to effectively utilize TensorCore WMMA capabilities. While LLMs could generate basic CUDA code, they often struggled to implement the specialized matrix multiplication acceleration features offered by modern GPU architectures. This suggests a potential gap in the training data or the models' understanding of advanced hardware-specific optimizations.
- Maintaining kernel proposal diversity proved challenging. The LLMs frequently gravitated toward similar optimization strategies, potentially missing more innovative solutions. This raises questions about optimal exploration-exploitation trade-offs in the kernel optimization process and suggests the need for more sophisticated diversity-promoting mechanisms.
- Practical constraints limited our parallel kernel verification to a population size of $N = 4$. While the per-GPU workload is relatively small, profiling operations often took longer than runtime estimation and compilation. This bottleneck in the evaluation pipeline suggests opportunities for more efficient verification methods.
- Selecting relevant feedback from profiling tools remains an open challenge. Determining which profiling metrics most effectively guide the optimization process requires further investigation.

**Broader Impact and Ethical Considerations**  Our work contributes to the ongoing trend of automating specialized engineering work. As highlighted by Handa et al. (2025), such automation can have significant economic and workforce implications. While tools like THE AI CUDA ENGINEER can enhance productivity and make specialized optimization more accessible, they may also impact the job market for CUDA optimization engineers. We also emphasize the critical importance of human verification for LLM-generated kernels. During our experiments, we occasionally discovered that LLMs could exploit non-robust correctness tests, highlighting the need for comprehensive validation procedures. Any deployment of automatically generated CUDA kernels should include thorough testing in sandboxed environments and careful human review, especially for safety-critical applications. These considerations suggest promising directions for future research while highlighting the importance of

responsible development and deployment of automated kernel optimization systems.

## 12. Discussion

**Costs & Runtime of THE AI CUDA ENGINEER**  Our results were produced with an estimated per-kernel cost of $15 in API credits for foundation model usage and less than 2 hours of runtime including kernel generation, verification, and profiling. Costs and runtime can be scaled to increase translation/optimization effectiveness as well as cover more complex tasks.

**Future Directions**  While this work focuses on optimizing the usage of frontier models, there is more work to be done in improving the underlying models at train time. With this project, we are releasing a large verified kernel library containing code, speedup times, profiling and error messages. Such a library could be used to train open-source models via supervised fine-tuning, preference optimization, or off-policy reinforcement learning (RL). Our verification and profiling setup can also be used to further train a model as in (Guo et al., 2025), which showed drastic gains in verifiable tasks through the simple scaling of online RL.

Additionally, there are several improvements to be made which are agnostic to the choice of train or inference time kernels. This includes hardware specialization – producing kernels that target the specific strengths or capabilities of the user device. Such specialization has become increasingly important for improving kernel efficiency. Flash Attention 2 (Dao, 2024), a prominent attention kernel, reaches 35% theoretical max FLOPS utilization on H100s, which are based on NVIDIA's Hopper architecture. Flash Attention 3 (Shah et al., 2024), which targets Hopper specifically, reaches 75% utilization and is 1.5-2x faster than Flash Attention 2 on Hopper devices.

Another exciting direction is the usage of compilers to aid the LLM in translating and optimizing modules into kernels. E.g., `torch` compile uses TorchDynamo to extract an FX graph from PyTorch code and jit-compile it to Triton (Ansel et al., 2024). Intuitively, these compilers translate PyTorch code to an intermediate representation that is more conducive to optimization and then undergo the actual optimization. It is possible that some or all of this work could be commandeered to aid THE AI CUDA ENGINEER in its tasks.

**Conclusion**  In this work, we introduced the AI CUDA Engineer, a comprehensive framework for automatic CUDA kernel discovery and optimization. Our approach demonstrates that large language models, when combined with evolutionary optimization strategies and retrieval-augmented generation, can successfully translate PyTorch operations into efficient CUDA implementations and further optimize their runtime performance. The AI CUDA Engineer achieves this through a four-stage pipeline that handles the conversion of PyTorch modules to functional representations, translation to CUDA kernels, runtime optimization, and kernel composition.

Looking forward, the AI CUDA Engineer opens up several promising research directions, including extending support for backward pass computations, hardware-specific optimizations, and deeper integration with existing compiler frameworks. While our current work focuses on inference-time optimizations, the principles and techniques developed here could be adapted to address the broader challenges of full-model optimization, including training-time computations and specialized architecture support.

As the demand for efficient deep learning computations continues to grow, tools like the AI CUDA Engineer demonstrate the potential for automated systems to assist in the complex task of performance optimization. This work represents a step toward more accessible and efficient hardware acceleration for deep learning workloads while highlighting the importance of combining traditional optimization techniques with modern language model capabilities.

## Acknowledgments

## References

Shubham Agarwal, Issam H Laradji, Laurent Charlin, and Christopher Pal. Litllm: A toolkit for scientific literature review. *arXiv preprint arXiv:2402.01788*, 2024.

Nurshat Fateh Ali, Md Mahdi Mohtasim, Shakil Mosharrof, and T Gopi Krishna. Automated literature review using nlp techniques and llm-based retrieval-augmented generation. *arXiv preprint arXiv:2411.18583*, 2024.

AlphaProof and AlphaGeometry teams. AI Achieves Silver-Medal Standard Solving International Mathematical Olympiad Problems. https://deepmind.google/discover/blog/ai-solves-imo-problems-at-silver-medal-level/, July 2024. Accessed: 2025-02-08.

Jason Ansel, Edward Yang, Horace He, Natalia Gimelshein, Animesh Jain, Michael Voznesensky, Bin Bao, Peter Bell, David Berard, Evgeni Burovski, Geeta Chauhan, Anjali Chourdia, Will Constable, Alban Desmaison, Zachary DeVito, Elias Ellison, Will Feng, Jiong Gong, Michael Gschwind, Brian Hirsh, Sherlock Huang, Kshiteej Kalambarkar, Laurent Kirsch, Michael Lazos, Mario Lezcano, Yanbo Liang, Jason Liang, Yinghai Lu, CK Luk, Bert Maher, Yunjie Pan, Christian Puhrsch, Matthias Reso, Mark Saroufim, Marcos Yukio Siraichi, Helen Suk, Michael Suo, Phil Tillet, Eikan Wang, Xiaodong Wang, William Wen, Shunting Zhang, Xu Zhao, Keren Zhou, Richard Zou, Ajit Mathews, Gregory Chanan, Peng Wu, and Soumith Chintala. PyTorch 2: Faster Machine Learning Through Dynamic Python Bytecode Transformation and Graph Compilation. In *Proceedings of the 29th ACM International Conference on Architectural Support for Programming Languages and Operating Systems, Volume 2 (ASPLOS '24)*. ACM, April 2024. doi: 10.1145/3620665.3640366. URL https://pytorch.org/assets/pytorch2-2.pdf.

Anthropic. Model card and evaluations for claude models, 2023. URL https://www-files.anthropic.com/production/images/Model-Card-Claude-2.pdf.

Anthropic. The claude 3 model family: Opus, sonnet, haiku, 2024. URL https://www-cdn.anthropic.com/de8ba9b01c9ab7cbabf5c33b80b7bbc618857627/Model_Card_Claude_3.pdf.

Jacob Austin, Augustus Odena, Maxwell Nye, Maarten Bosma, Henryk Michalewski, David Dohan, Ellen Jiang, Carrie Cai, Michael Terry, Quoc Le, et al. Program synthesis with large language models. *arXiv preprint arXiv:2108.07732*, 2021.

Jeremy Berman. How i got a record 53.6% on arc-agi. https://jeremyberman.substack.com/p/how-i-got-a-record-536-on-arc-agi, 2025. Accessed: 2025-02-08.

Rishi Bommasani, Drew A Hudson, Ehsan Adeli, Russ Altman, Simran Arora, Sydney von Arx, Michael S Bernstein, Jeannette Bohg, Antoine Bosselut, Emma Brunskill, et al. On the opportunities and risks of foundation models. *arXiv preprint arXiv:2108.07258*, 2021.

Bradley Brown, Jordan Juravsky, Ryan Ehrlich, Ronald Clark, Quoc V Le, Christopher Ré, and Azalia Mirhoseini. Large language monkeys: Scaling inference compute with repeated sampling. *arXiv preprint arXiv:2407.21787*, 2024.

Terry Chen, Bing Xu, and Kirthi Devleker. Automating gpu kernel generation with deepseek-r1 and inference time scaling, February 2025. URL https://developer.nvidia.com/blog/automating-gpu-kernel-generation-with-deepseek-r1-and-inference-time-scaling/?ncid=so-twit-997075&linkId=100000338909937.

Tianqi Chen, Thierry Moreau, Ziheng Jiang, Lianmin Zheng, Eddie Yan, Haichen Shen, Meghan Cowan, Leyuan Wang, Yuwei Hu, Luis Ceze, et al. {TVM}: An automated {End-to-End} optimizing compiler for deep learning. In *13th USENIX Symposium on Operating Systems Design and Implementation (OSDI 18)*, pages 578–594, 2018.

Sharan Chetlur, Cliff Woolley, Philippe Vandermersch, Jonathan Cohen, John Tran, Bryan Catanzaro, and Evan Shelhamer. cudnn: Efficient primitives for deep learning. *arXiv preprint arXiv:1410.0759*, 2014.

Tri Dao. Flashattention-2: Faster attention with better parallelism and work partitioning. In *The Twelfth International Conference on Learning Representations*, 2024. URL https://openreview.net/forum?id=mZn2Xyh9Ec.

Debajyoti Dasgupta, Arijit Mondal, and Partha Pratim Chakrabarti. Empowering AI as autonomous researchers: Evaluating LLMs in generating novel research ideas through automated metrics. In *2nd AI4Research Workshop: Towards a Knowledge-grounded Scientific Research Lifecycle*, 2024. URL https://openreview.net/forum?id=12T3Nt22av.

Yilun Du, Shuang Li, Antonio Torralba, Joshua B Tenenbaum, and Igor Mordatch. Improving factuality and reasoning in language models through multiagent debate. *arXiv preprint arXiv:2305.14325*, 2023.

Ryan Ehrlich, Bradley Brown, Jordan Juravsky, Ronald Clark, Christopher Ré, and Azalia Mirhoseini. Codemonkeys: Scaling test-time compute for software engineering. *arXiv preprint arXiv:2501.14723*, 2025.

Martin Ester, Hans-Peter Kriegel, Jörg Sander, Xiaowei Xu, et al. A density-based algorithm for discovering clusters in large spatial databases with noise. In *kdd*, volume 96, pages 226–231, 1996.

Maxence Faldor, Jenny Zhang, Antoine Cully, and Jeff Clune. Omni-epic: Open-endedness via models of human notions of interestingness with environments programmed in code, 2024. URL https://arxiv.org/abs/2405.15568.

Yunfan Gao, Yun Xiong, Xinyu Gao, Kangxiang Jia, Jinliu Pan, Yuxi Bi, Yi Dai, Jiawei Sun, and Haofen Wang. Retrieval-augmented generation for large language models: A survey. *arXiv preprint arXiv:2312.10997*, 2023.

Michael Garland, Scott Le Grand, John Nickolls, Joshua Anderson, Jim Hardwick, Scott Morton, Everett Phillips, Yao Zhang, and Vasily Volkov. Parallel computing experiences with cuda. *IEEE micro*, 28(4):13–27, 2008.

Paul Gauthier. aider, 2024. URL https://github.com/paul-gauthier/aider.

Daya Guo, Dejian Yang, Haowei Zhang, Junxiao Song, Ruoyu Zhang, Runxin Xu, Qihao Zhu, Shirong Ma, Peiyi Wang, Xiao Bi, et al. Deepseek-r1: Incentivizing reasoning capability in llms via reinforcement learning. *arXiv preprint arXiv:2501.12948*, 2025.

Kunal Handa, Alex Tamkin, Miles McCain, Saffron Huang, Esin Durmus, Sarah Heck, Jared Mueller, Jerry Hong, Stuart Ritchie, Tim Belonax, Kevin K Troy, Dario Amodei, Jared Kaplan, Jack Clark,

and Deep Ganguli. Which economic tasks are performed with AI? evidence from millions of claude conversations. 2025.

Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 770–778, 2016.

Danny Hernandez and Tom B Brown. Measuring the algorithmic efficiency of neural networks. *arXiv preprint arXiv:2005.04305*, 2020.

Xun Huang and Serge Belongie. Arbitrary style transfer in real-time with adaptive instance normalization. In *Proceedings of the IEEE international conference on computer vision*, pages 1501–1510, 2017.

Binyuan Hui, Jian Yang, Zeyu Cui, Jiaxi Yang, Dayiheng Liu, Lei Zhang, Tianyu Liu, Jiajun Zhang, Bowen Yu, Keming Lu, et al. Qwen2. 5-coder technical report. *arXiv preprint arXiv:2409.12186*, 2024.

Aaron Jaech, Adam Kalai, Adam Lerer, Adam Richardson, Ahmed El-Kishky, Aiden Low, Alec Helyar, Aleksander Madry, Alex Beutel, Alex Carney, et al. Openai o1 system card. *arXiv preprint arXiv:2412.16720*, 2024.

Carlos E Jimenez, John Yang, Alexander Wettig, Shunyu Yao, Kexin Pei, Ofir Press, and Karthik R Narasimhan. SWE-bench: Can language models resolve real-world github issues? In *The Twelfth International Conference on Learning Representations*, 2024. URL https://openreview.net/forum?id=VTF8yNQM66.

Norman P Jouppi, Cliff Young, Nishant Patil, David Patterson, Gaurav Agrawal, Raminder Bajwa, Sarah Bates, Suresh Bhatia, Nan Boden, Al Borchers, et al. In-datacenter performance analysis of a tensor processing unit. In *Proceedings of the 44th annual international symposium on computer architecture*, pages 1–12, 2017.

Jared Kaplan, Sam McCandlish, Tom Henighan, Tom B Brown, Benjamin Chess, Rewon Child, Scott Gray, Alec Radford, Jeffrey Wu, and Dario Amodei. Scaling laws for neural language models. *arXiv preprint arXiv:2001.08361*, 2020.

David B Kirk and W Hwu Wen-Mei. *Programming massively parallel processors: a hands-on approach*. Morgan kaufmann, 2016.

Robert Lange, Tom Schaul, Yutian Chen, Tom Zahavy, Valentin Dalibard, Chris Lu, Satinder Singh, and Sebastian Flennerhag. Discovering evolution strategies via meta-black-box optimization. In *Proceedings of the Companion Conference on Genetic and Evolutionary Computation*, pages 29–30, 2023.

Robert Tjarko Lange, Yingtao Tian, and Yujin Tang. Large language models as evolution strategies. *arXiv preprint arXiv:2402.18381*, 2024.

Joel Lehman, Jonathan Gordon, Shawn Jain, Kamal Ndousse, Cathy Yeh, and Kenneth O. Stanley. Evolution through large models, 2022. URL https://arxiv.org/abs/2206.08896.

Patrick Lewis, Ethan Perez, Aleksandra Piktus, Fabio Petroni, Vladimir Karpukhin, Naman Goyal, Heinrich Küttler, Mike Lewis, Wen-tau Yih, Tim Rocktäschel, Sebastian Riedel, and Douwe Kiela. Retrieval-augmented generation for knowledge-intensive nlp tasks. In *Proceedings of the 34th International Conference on Neural Information Processing Systems*, NIPS '20, Red Hook, NY, USA, 2020a. Curran Associates Inc. ISBN 9781713829546.

Patrick Lewis, Ethan Perez, Aleksandra Piktus, Fabio Petroni, Vladimir Karpukhin, Naman Goyal, Heinrich Küttler, Mike Lewis, Wen-tau Yih, Tim Rocktäschel, et al. Retrieval-augmented generation for knowledge-intensive nlp tasks. *Advances in Neural Information Processing Systems*, 33:9459–9474, 2020b.

Wen-Ding Li, Keya Hu, Carter Larsen, Yuqing Wu, Simon Alford, Caleb Woo, Spencer M. Dunn, Hao Tang, Wei-Long Zheng, Yewen Pu, and Kevin Ellis. Combining induction and transduction for abstract reasoning. In *The Thirteenth International Conference on Learning Representations*, 2025. URL https://openreview.net/forum?id=UmdotAAVDe.

Yujia Li, David Choi, Junyoung Chung, Nate Kushman, Julian Schrittwieser, Rémi Leblond, Tom Eccles, James Keeling, Felix Gimeno, Agustin Dal Lago, et al. Competition-level code generation with alphacode. *Science*, 378(6624):1092–1097, 2022.

Aixin Liu, Bei Feng, Bing Xue, Bingxuan Wang, Bochao Wu, Chengda Lu, Chenggang Zhao, Chengqi Deng, Chenyu Zhang, Chong Ruan, et al. Deepseek-v3 technical report. *arXiv preprint arXiv:2412.19437*, 2024.

Ryan Liu and Nihar Shah. Reviewergpt? an exploratory study on using large language models for paper reviewing. *arXiv preprint arXiv:2306.00622*, 2023.

Chris Lu, Samuel Holt, Claudio Fanconi, Alex James Chan, Jakob Nicolaus Foerster, Mihaela van der Schaar, and Robert Tjarko Lange. Discovering preference optimization algorithms with and for large language models. In *The Thirty-eighth Annual Conference on Neural Information Processing Systems*, 2024a. URL https://openreview.net/forum?id=erjQDJOz9L.

Chris Lu, Cong Lu, Robert Tjarko Lange, Jakob Foerster, Jeff Clune, and David Ha. The ai scientist: Towards fully automated open-ended scientific discovery. *arXiv preprint arXiv:2408.06292*, 2024b.

Cong Lu, Shengran Hu, and Jeff Clune. Beyond benchmarking: Automated capability discovery via model self-exploration. In *Language Gamification-NeurIPS 2024 Workshop*.

David Luebke. Cuda: Scalable parallel programming for high-performance scientific computing. In *2008 5th IEEE international symposium on biomedical imaging: from nano to macro*, pages 836–838. IEEE, 2008.

Liangchen Luo, Yinxiao Liu, Rosanne Liu, Samrat Phatale, Harsh Lara, Yunxuan Li, Lei Shu, Yun Zhu, Lei Meng, Jiao Sun, et al. Improve mathematical reasoning in language models by automated process supervision. *arXiv preprint arXiv:2406.06592*, 2024.

Yecheng Jason Ma, William Liang, Guanzhi Wang, De-An Huang, Osbert Bastani, Dinesh Jayaraman, Yuke Zhu, Linxi Fan, and Anima Anandkumar. Eureka: Human-level reward design via coding large language models. *arXiv preprint arXiv:2310.12931*, 2023.

METR. Measuring automated kernel engineering, February 2025. URL https://metr.org/blog/2025-02-14-measuring-automated-kernel-engineering/.

Elliot Meyerson, Mark J Nelson, Herbie Bradley, Adam Gaier, Arash Moradi, Amy K Hoover, and Joel Lehman. Language model crossover: Variation through few-shot prompting. *arXiv preprint arXiv:2302.12170*, 2023.

Niklas Muennighoff, Zitong Yang, Weijia Shi, Xiang Lisa Li, Li Fei-Fei, Hannaneh Hajishirzi, Luke Zettlemoyer, Percy Liang, Emmanuel Candès, and Tatsunori Hashimoto. s1: Simple test-time scaling. *arXiv preprint arXiv:2501.19393*, 2025.

John Nickolls, Ian Buck, Michael Garland, and Kevin Skadron. Scalable parallel programming with cuda: Is cuda the parallel programming model that application developers have been waiting for? *Queue*, 6(2):40–53, 2008.

Anne Ouyang, Simon Guo, Simran Arora, Alex L. Zhang, William Hu, Christopher Ré, and Azalia Mirhoseini. Kernelbench: Can llms write efficient gpu kernels?, 2025. URL https://arxiv.org/abs/2502.10517.

Shanghaoran Quan, Jiaxi Yang, Bowen Yu, Bo Zheng, Dayiheng Liu, An Yang, Xuancheng Ren, Bofei Gao, Yibo Miao, Yunlong Feng, et al. Codeelo: Benchmarking competition-level code generation of llms with human-comparable elo ratings. *arXiv preprint arXiv:2501.01257*, 2025.

Matthew Renze and Erhan Guven. The effect of sampling temperature on problem solving in large language models. *arXiv preprint arXiv:2402.05201*, 2024.

Bernardino Romera-Paredes, Mohammadamin Barekatain, Alexander Novikov, Matej Balog, M Pawan Kumar, Emilien Dupont, Francisco JR Ruiz, Jordan S Ellenberg, Pengming Wang, Omar Fawzi, et al. Mathematical discoveries from program search with large language models. *Nature*, 625(7995): 468–475, 2024.

Baptiste Roziere, Jonas Gehring, Fabian Gloeckle, Sten Sootla, Itai Gat, Xiaoqing Ellen Tan, Yossi Adi, Jingyu Liu, Romain Sauvestre, Tal Remez, et al. Code llama: Open foundation models for code. *arXiv preprint arXiv:2308.12950*, 2023.

Jaime Sevilla, Lennart Heim, Anson Ho, Tamay Besiroglu, Marius Hobbhahn, and Pablo Villalobos. Compute trends across three eras of machine learning. In *2022 International Joint Conference on Neural Networks (IJCNN)*, pages 1–8. IEEE, 2022.

Jay Shah, Ganesh Bikshandi, Ying Zhang, Vijay Thakkar, Pradeep Ramani, and Tri Dao. Flashattention-3: Fast and accurate attention with asynchrony and low-precision. In *The Thirty-eighth Annual Conference on Neural Information Processing Systems*, 2024. URL https://openreview.net/forum?id=tVConYid20.

Chenglei Si, Diyi Yang, and Tatsunori Hashimoto. Can llms generate novel research ideas? *arXiv preprint arXiv:2409.04109*, 2024.

Charlie Snell, Jaehoon Lee, Kelvin Xu, and Aviral Kumar. Scaling llm test-time compute optimally can be more effective than scaling model parameters. *arXiv preprint arXiv:2408.03314*, 2024.

Xingyou Song, Yingtao Tian, Robert Tjarko Lange, Chansoo Lee, Yujin Tang, and Yutian Chen. Position paper: Leveraging foundational models for black-box optimization: Benefits, challenges, and future directions. *arXiv preprint arXiv:2405.03547*, 2024.

Benjamin F Spector, Simran Arora, Aaryan Singhal, Daniel Y Fu, and Christopher Ré. Thunderkittens: Simple, fast, and adorable ai kernels. *arXiv preprint arXiv:2410.20399*, 2024.

Kenneth O Stanley and Joel Lehman. *Why greatness cannot be planned: The myth of the objective*. Springer, 2015.

Teo Susnjak, Peter Hwang, Napoleon H. Reyes, Andre L. C. Barczak, Timothy R. McIntosh, and Surangika Ranathunga. Automating research synthesis with domain-specific large language model fine-tuning. *ACM Trans. Knowl. Discov. Data*, January 2025. ISSN 1556-4681. doi: 10.1145/3715 964. URL https://doi.org/10.1145/3715964. Just Accepted.

Kimi Team, Angang Du, Bofei Gao, Bowei Xing, Changjiu Jiang, Cheng Chen, Cheng Li, Chenjun Xiao, Chenzhuang Du, Chonghua Liao, et al. Kimi k1. 5: Scaling reinforcement learning with llms. *arXiv preprint arXiv:2501.12599*, 2025.

Philippe Tillet, Hsiang-Tsung Kung, and David Cox. Triton: an intermediate language and compiler for tiled neural network computations. In *Proceedings of the 3rd ACM SIGPLAN International Workshop on Machine Learning and Programming Languages*, pages 10–19, 2019.

Fouad Trad and Ali Chehab. To ensemble or not: Assessing majority voting strategies for phishing detection with large language models. *arXiv preprint arXiv:2412.00166*, 2024.

Laurens Van der Maaten and Geoffrey Hinton. Visualizing data using t-sne. *Journal of machine learning research*, 9(11), 2008.

Junlin Wang, Jue WANG, Ben Athiwaratkun, Ce Zhang, and James Zou. Mixture-of-agents enhances large language model capabilities. In *The Thirteenth International Conference on Learning Representations*, 2025. URL https://openreview.net/forum?id=h0ZfDIrj7T.

Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten Bosma, Fei Xia, Ed Chi, Quoc V Le, Denny Zhou, et al. Chain-of-thought prompting elicits reasoning in large language models. *Advances in neural information processing systems*, 35:24824–24837, 2022.

Yuxi Xie, Kenji Kawaguchi, Yiran Zhao, Xu Zhao, Min-Yen Kan, Junxian He, and Qizhe Xie. Self-evaluation guided beam search for reasoning. In *Thirty-seventh Conference on Neural Information Processing Systems*, 2023. URL https://openreview.net/forum?id=Bw82hwg5Q3.

Zonglin Yang, Xinya Du, Junxian Li, Jie Zheng, Soujanya Poria, and Erik Cambria. Large language models for automated open-domain scientific hypotheses discovery. *arXiv preprint arXiv:2309.02726*, 2023.

Dan Zhang, Sining Zhoubian, Ziniu Hu, Yisong Yue, Yuxiao Dong, and Jie Tang. ReST-MCTS*: LLM self-training via process reward guided tree search. In *The Thirty-eighth Annual Conference on Neural Information Processing Systems*, 2024a. URL https://openreview.net/forum?id=8rcFOqEud5.

Genghan Zhang, Weixin Liang, Olivia Hsu, and Kunle Olukotun. Adaptive self-improvement llm agentic system for ml library development. *arXiv preprint arXiv:2502.02534*, 2025.

Jenny Zhang, Joel Lehman, Kenneth Stanley, and Jeff Clune. OMNI: Open-endedness via models of human notions of interestingness. In *The Twelfth International Conference on Learning Representations*, 2024b. URL https://openreview.net/forum?id=AgM3MzT99c.

Lianmin Zheng, Chengfan Jia, Minmin Sun, Zhao Wu, Cody Hao Yu, Ameer Haj-Ali, Yida Wang, Jun Yang, Danyang Zhuo, Koushik Sen, et al. Ansor: Generating {High-Performance} tensor programs for deep learning. In *14th USENIX symposium on operating systems design and implementation (OSDI 20)*, pages 863–879, 2020.

Denny Zhou, Nathanael Schärli, Le Hou, Jason Wei, Nathan Scales, Xuezhi Wang, Dale Schuurmans, Claire Cui, Olivier Bousquet, Quoc Le, et al. Least-to-most prompting enables complex reasoning in large language models. *arXiv preprint arXiv:2205.10625*, 2022.

Yifei Zhou, Andrea Zanette, Jiayi Pan, Sergey Levine, and Aviral Kumar. ArCHer: Training language model agents via hierarchical multi-turn RL. In *Forty-first International Conference on Machine Learning*, 2024. URL https://openreview.net/forum?id=b6rA0kAHT1.

Qihao Zhu, Daya Guo, Zhihong Shao, Dejian Yang, Peiyi Wang, Runxin Xu, Y Wu, Yukun Li, Huazuo Gao, Shirong Ma, et al. Deepseek-coder-v2: Breaking the barrier of closed-source models in code intelligence. *arXiv preprint arXiv:2406.11931*, 2024.

Zhenzhen Zhuang, Jiandong Chen, Hongfeng Xu, Yuwen Jiang, and Jialiang Lin. Large language models for automated scholarly paper review: A survey. *arXiv preprint arXiv:2501.10326*, 2025.

# Appendix

# A. Additional Results

## A.1. Level 1: Improvement Summary



Figure 19 | **Level 1 Full Results.**

## A.2. Level 1: Evolutionary Optimization Learning Curves



Figure 20 | **Level 1 Optimization Curves.**

Figure 21 | **Level 1 Optimization Curves.**

## A.3. Level 2: Improvement Summary



Figure 22 | **Level 2 Full Results.**

## A.4. Level 2: Evolutionary Optimization Learning Curves



Figure 23 | **Level 2 Optimization Curves.**

Figure 24 | **Level 2 Optimization Curves.**

## A.5. Level 3: Improvement Summary



Figure 25 | **Level 3 Full Results.**

## A.6. Level 3: Evolutionary Optimization Learning Curves



Figure 26 | **Level 3 Optimization Curves.**

# B. Prompts

We present some representative prompts that we use for THE AI CUDA ENGINEER.

## B.1. Conversion Prompts

These prompts correspond to the first stage of THE AI CUDA ENGINEER.

---

**Conversion System Prompt**

```
You are an expert python and PyTorch engineer.

Your job is correctness and holding to the given task specification.
```

---

**Conversion Base Prompt**

```
You will be given python code, which contains four parts:
- [Imports] These are always at the top of the file and contains
↪   mostly torch imports.
- [Model Definition] A `Model(nn.Module)` class with an init and a
↪   forward method.
- [Configurations] Following the `Model` class, there are some
↪   configurations and two functions `get_inputs()` and
↪   `get_init_inputs()`. The configurations are used in these
↪   functions.
- [Members] The `Model` class may have members that are also
↪   `nn.Module`, they are defined either in `torch.nn` or before the
↪   `Model` class.

Your job is to write a "functional" version of that code, which
↪   suggests the workflow of
↪   `Model(*get_init_inputs())(*get_inputs())`.
The task for each part of the code is as follows:

For [Imports], you will likely need the following libraries:
```python
import torch
import torch.nn as nn
import torch.nn.functional as F
```
You may import other libraries, but you probably not any other torch
↪   modules.
For example, if you do `from torch import _VF`, you can then use
↪   `_VF.lstm` and `_VF.gru` functions.

For [Model Definition], you should:
- Define `nn.Parameter` in `Model.__init__()`, since they are needed
↪   for the functional calls.
    - If a parameter is already defined in the `Model` class, you can
    ↪   directly use it.
```

---

```
        - Otherwise, you may extract it from the `nn.Module`s. E.g.
        ↪  `self.conv.weights =
        ↪  nn.Parameter(self.conv1.weight.data.clone)`, where
        ↪  `self.conv1` is an `nn.Module`. If you use `nn.ParameterDict`
        ↪  to extract parameters, remember that names cannot contain
        ↪  dots.
        - Notice that these parameters' specifications are given by the
        ↪  `get_init_inputs()` function call because Model is
        ↪  initialized with `Model(*get_init_inputs())`.
- Modify the `forward()` method to use functional calls instead of
↪  `nn.Modules`.
        - This can be achieved by defining a `module_fn()` that takes in
        ↪  all the args that the forward pass was called with as well as
        ↪  any neural network parameters that the Model holds. It should
        ↪  faithfully reproduce the exact forward pass as the original
        ↪  Model class.
        - Then, you can augment the `Model.forward()` method signature,
        ↪  where you add a `fn` argument that defaults to `module_fn`.
- You need to read the original code carefully. E.g., the code may
↪  not implement a residual connection even though the name or
↪  comment suggests it.


For [Configurations], you should simply copy the code from the
↪  original file. They are for test purposes.


For [Members], you should:
- If a member is an `nn.Module` for which there exists a functional
↪  call, you can replace its forward pass with the functional call.
↪  E.g., if it is a `nn.Conv2d`, you can replace `self.conv(x)` with
↪  `F.conv2d(x, ...)` with proper arguments.
- Otherwise, you need to decompose it until the its forward pass can
↪  be replaced with a series of functional calls. You can define
↪  functions for this, and these functions will be used in the
↪  `module_fn()` function.


Your returned functional version of the code should be a valid python
↪  file, and it will be checked against the original code. Their
↪  outputs should be identical.
Return only python code, no other text.
The following are several examples to illustrate the task.


=== Example 1 ===


{example1}


=== Example 2 ===


{example2}
```

```
=== Example 3 ===

{example3}

=== Example 4 ===

{example4}

=================

Here is the code you need to convert:
```python
{code}
```
```

## B.2. Translation Prompts

These prompts correspond to the second stage of THE AI CUDA ENGINEER.

> **Translation System Prompt**
>
> ```
> You are a CUDA engineer tasked with translating PyTorch code into
> ↪  CUDA kernel code.
>
> The CUDA code you generate will be saved in `cuda_fname` and loaded
> ↪  using torch.utils.cpp_extension.load():
> ```python
> cuda_fn = load(
>     name=task_name,
>     sources=[cuda_fname],
>     extra_cuda_cflags=["-O3", "--use_fast_math"],
>     with_cuda=True,
>     verbose=True,
> )
> ```
> Later, the function will be called via `cuda_fn =
> ↪  load(name=task_name, ...).forward` and thoroughly tested.
>
> Translate the PyTorch code (in <pytorch> tags) into CUDA kernel code.
>
> <instructions>
> - Write CUDA code that performs the **exact same operation** as the
> ↪  PyTorch code.
> - Include the required pybind11 cuda module name in the code.
> - Return the code between <cuda></cuda> tags.
> </instructions>
> ```

**Translation Base Prompt**

```
<pytorch>
{EXAMPLES}
</pytorch>

Now, translate the PyTorch code into CUDA kernel code.
```

## B.3. Optimization Prompts

**Optimization System Prompt**

```
You are a Machine Learning Engineer trying to reduce the runtime of a
↪ {operation} kernel in CUDA. Make sure the kernel returns the
↪  correct result. Do not use any alternative precision that could
↪  result in an incorrect result. The kernel will be run on a
↪ {gpu_type} GPU with CUDA {cuda_version}. {mode_prompt}

Answer using the following schema:

name: A shortened descriptor of the idea. Lowercase, no spaces,
↪  underscores allowed.
code: The proposed cuda script in code.
thought: The rationale for the improvement idea.

The pybind11 cuda module name has to be the same as in the example.
MAKE SURE THE PROPOSAL CODE IS VALID CUDA CODE.
FOLLOW EXACTLY THIS FORMAT. DO NOT ADD ANYTHING ELSE.

Here is a reference implementation of a working CUDA kernel proposal:

name: {name_example}
thought: {thought_example}
code: {code_example}

runtime: {runtime_example:0.5f} milliseconds
speedup over torch: {speedup_example:0.2f}x
```

**Optimization Base Prompt**

```
Propose a new CUDA kernel (including name, code, thought) which aims
↪  to reduce the runtime of the operation, while ensuring the kernel
↪  returns the correct result.
```

# C. Highlighted Discovered Kernels - KernelBench Level 1

## C.1. 34_instance_norm

```
#include <torch/extension.h>
#include <cuda.h>
#include <cuda_runtime.h>
#include <vector>

// Warp-level reduction using shuffle instructions
__inline__ __device__ float warpReduceSum(float val) {
    for (int offset = warpSize / 2; offset > 0; offset /= 2) {
        val += __shfl_down_sync(0xffffffff, val, offset);
    }
    return val;
}

// Block-level reduction using a statically allocated shared memory array for reduction
__inline__ __device__ float blockReduceSum(float val) {
    __shared__ float s_reduce[32];
    int lane = threadIdx.x % warpSize;
    int wid  = threadIdx.x / warpSize;

    val = warpReduceSum(val);
    if (lane == 0) {
        s_reduce[wid] = val;
    }
    __syncthreads();

    val = (threadIdx.x < blockDim.x / warpSize) ? s_reduce[lane] : 0.0f;
    if (wid == 0) {
        val = warpReduceSum(val);
    }
    return val;
}

// Kernel that leverages shared memory to stage the instance data, reducing global memory accesses
__global__ void instance_norm_kernel_shared(
    const float* __restrict__ x,
    float* __restrict__ y,
    const float* __restrict__ weight,
    const float* __restrict__ bias,
    const int N,
    const int C,
    const int H,
    const int W,
    const float eps
) {
    const int HW = H * W;
    const int instance_idx = blockIdx.x; // one block per (N, C) instance
    if (instance_idx >= N * C) return;

    int n = instance_idx / C;
    int c = instance_idx % C;
    const int instance_offset = (n * C + c) * HW;

    // Load scale and shift from weight and bias using __ldg for caching
    const float scale = weight ? __ldg(&weight[c]) : 1.0f;
    const float shift = bias ? __ldg(&bias[c]) : 0.0f;

    // Allocate shared memory for the instance data
    extern __shared__ float s_data[];  // size = HW * sizeof(float)
    float* s_x = s_data;

    // Load the instance data from global memory into shared memory
    for (int i = threadIdx.x; i < HW; i += blockDim.x) {
        s_x[i] = x[instance_offset + i];
    }
    __syncthreads();

    // Compute local sums and squared sums from shared memory
    float local_sum = 0.0f;
    float local_sq_sum = 0.0f;
    for (int i = threadIdx.x; i < HW; i += blockDim.x) {
        float val = s_x[i];
        local_sum += val;
        local_sq_sum += val * val;
    }

    // Block-level reduction to compute sum and squared sum
    float sum = blockReduceSum(local_sum);
    float sum_sq = blockReduceSum(local_sq_sum);

    __shared__ float s_mean;
    __shared__ float s_inv_std;
    if (threadIdx.x == 0) {
        float mean = sum / HW;
        float var = fmaxf(sum_sq / HW - mean * mean, 0.0f);
        s_mean = mean;
        s_inv_std = rsqrtf(var + eps);
    }
    __syncthreads();

    // Normalize the data in shared memory using the computed mean and inverse std
    for (int i = threadIdx.x; i < HW; i += blockDim.x) {
```

```
            float val = s_x[i];
            s_x[i] = ((val - s_mean) * s_inv_std * scale) + shift;
        }
        __syncthreads();

        // Write the normalized data back to global memory
        for (int i = threadIdx.x; i < HW; i += blockDim.x) {
            y[instance_offset + i] = s_x[i];
        }
}

// Forward function called from PyTorch
torch::Tensor forward(
    torch::Tensor x,
    torch::Tensor weight,
    torch::Tensor bias,
    double eps
) {
    TORCH_CHECK(x.is_cuda(), "x must be a CUDA tensor");
    if (weight.defined()) TORCH_CHECK(weight.is_cuda(), "weight must be a CUDA tensor");
    if (bias.defined()) TORCH_CHECK(bias.is_cuda(), "bias must be a CUDA tensor");

    auto sizes = x.sizes();
    TORCH_CHECK(sizes.size() == 4, "Input tensor must be 4D: (N, C, H, W)");

    int N = sizes[0];
    int C = sizes[1];
    int H = sizes[2];
    int W = sizes[3];

    auto y = torch::empty_like(x);

    // One thread block per instance
    int blocks = N * C;
    int threads = 256;
    // Allocate shared memory: one instance per block
    int shared_mem_size = H * W * sizeof(float);

    instance_norm_kernel_shared<<<blocks, threads, shared_mem_size>>>(
        x.data_ptr<float>(),
        y.data_ptr<float>(),
        weight.defined() ? weight.data_ptr<float>() : nullptr,
        bias.defined() ? bias.data_ptr<float>() : nullptr,
        N, C, H, W,
        static_cast<float>(eps)
    );

    return y;
}

PYBIND11_MODULE(TORCH_EXTENSION_NAME, m) {
    m.def("forward", &forward, "Instance Normalization forward (CUDA) with shared memory");
}
```

Listing 2 | THE AI CUDA ENGINEER Optimized CUDA Kernel for Instance Normalization.

## C.2. 15_matmul_lower

```
#include <torch/extension.h>
#include <cuda.h>
#include <cuda_runtime.h>

__global__ void triangular_mm_kernel(const float* __restrict__ A,
                                     const float* __restrict__ B,
                                     float* __restrict__ C,
                                     const int N) {
    // Use 2D block configuration for better occupancy
    const int row = blockIdx.y * blockDim.y + threadIdx.y;
    const int col = blockIdx.x * blockDim.x + threadIdx.x;

    if (row < N && col < N) {
        if (col <= row) {
            // Lower triangle computation
            float sum = 0.0f;

            // Process elements in chunks to improve cache utilization
            #pragma unroll 8
            for (int k = col; k <= row; k++) {
                sum += A[row * N + k] * B[k * N + col];
            }
            C[row * N + col] = sum;
        } else {
            // Upper triangle (set to zero)
            C[row * N + col] = 0.0f;
        }
    }
}

at::Tensor forward(at::Tensor A, at::Tensor B) {
    TORCH_CHECK(A.is_cuda(), "A must be a CUDA tensor");
    TORCH_CHECK(B.is_cuda(), "B must be a CUDA tensor");
    TORCH_CHECK(A.dim() == 2, "A must be a 2D tensor");
```

```
    TORCH_CHECK(B.dim() == 2, "B must be a 2D tensor");
    TORCH_CHECK(A.size(0) == A.size(1), "A must be square");
    TORCH_CHECK(B.size(0) == B.size(1), "B must be square");
    TORCH_CHECK(A.size(0) == B.size(0), "A and B must be the same size");

    int N = A.size(0);
    auto C = torch::empty_like(A);

    // Optimize thread count based on matrix size
    const int threadsPerBlock = 256;  // Increased thread count per block
    const int numBlocks = N;

    triangular_mm_kernel<<<numBlocks, threadsPerBlock>>>(
        A.data_ptr<float>(),
        B.data_ptr<float>(),
        C.data_ptr<float>(),
        N
    );

    cudaError_t err = cudaGetLastError();
    TORCH_CHECK(err == cudaSuccess, "CUDA kernel failed: ", cudaGetErrorString(err));

    return C;
}

PYBIND11_MODULE(TORCH_EXTENSION_NAME, m) {
    m.def("forward", &forward, "Strided efficient triangular matrix multiplication (CUDA)");
}
```

Listing 3 | THE AI CUDA ENGINEER Optimized CUDA Kernel for Lower Triangular Matrix Multiplication.

## C.3. 12_matmul_diagonal

```
#include <torch/extension.h>
#include <cuda.h>
#include <cuda_runtime.h>

__global__ void hybrid_diag_matmul_kernel(
    const float* __restrict__ A,
    const float* __restrict__ B,
    float* __restrict__ C,
    const int64_t N,
    const int64_t M,
    const bool use_vectorized
) {
    if (use_vectorized) {
        // Vectorized approach for large matrices where M is divisible by 4
        const int tid = threadIdx.x + blockIdx.x * blockDim.x;
        const int stride = blockDim.x * gridDim.x;
        const int total = N * M;
        const int vec_total = total / 4;

        const float4* B_vec = reinterpret_cast<const float4*>(B);
        float4* C_vec = reinterpret_cast<float4*>(C);

        for (int idx = tid; idx < vec_total; idx += stride) {
            const int base_idx = idx * 4;
            const int row = base_idx / M;
            const float a_val = A[row];

            float4 b_val = B_vec[idx];
            float4 c_val;
            c_val.x = a_val * b_val.x;
            c_val.y = a_val * b_val.y;
            c_val.z = a_val * b_val.z;
            c_val.w = a_val * b_val.w;

            C_vec[idx] = c_val;
        }
    } else {
        // Row-based approach for smaller matrices or when M is not divisible by 4
        int row = blockIdx.x;
        if (row < N) {
            float a_val = A[row];
            const int main_end = (M / blockDim.x) * blockDim.x;

            // Main loop with coalesced access
            for (int j = threadIdx.x; j < main_end; j += blockDim.x) {
                int idx = row * M + j;
                C[idx] = a_val * B[idx];
            }

            // Handle remaining elements
            for (int j = main_end + threadIdx.x; j < M; j += blockDim.x) {
                int idx = row * M + j;
                C[idx] = a_val * B[idx];
            }
        }
    }
}
```

```
at::Tensor forward(at::Tensor A, at::Tensor B) {
    TORCH_CHECK(A.dim() == 1, "A must be a 1D tensor");
    TORCH_CHECK(B.dim() == 2, "B must be a 2D tensor");
    TORCH_CHECK(A.size(0) == B.size(0), "Dimension mismatch");

    A = A.contiguous();
    B = B.contiguous();

    int64_t N = A.size(0);
    int64_t M = B.size(1);
    auto C = torch::empty({N, M}, B.options());

    // Choose approach based on matrix size and alignment
    bool use_vectorized = (M >= 512) && (M % 4 == 0);

    if (use_vectorized) {
        const int threads = 256;
        const int blocks = min(65535, (int)((N * M + threads * 4 - 1) / (threads * 4)));
        hybrid_diag_matmul_kernel<<<blocks, threads>>>(
            A.data_ptr<float>(), B.data_ptr<float>(), C.data_ptr<float>(),
            N, M, true);
    } else {
        int threads = (M > 256) ? 256 : (((M + 31) / 32) * 32);
        dim3 grid(N);
        hybrid_diag_matmul_kernel<<<grid, threads>>>(
            A.data_ptr<float>(), B.data_ptr<float>(), C.data_ptr<float>(),
            N, M, false);
    }

    return C;
}

PYBIND11_MODULE(TORCH_EXTENSION_NAME, m) {
    m.def("forward", &forward, "Hybrid diagonal matrix multiplication");
}
```

Listing 4 | THE AI CUDA ENGINEER Optimized CUDA Kernel for Diagonal Matrix Multiplication.

## C.4. 95_crossentropy

```
#include <torch/extension.h>

__global__ void cross_entropy_loss_kernel(
    const float* logits,
    const int64_t* targets,
    float* losses,
    int batch_size,
    int num_classes
)
{
    int i = blockIdx.x * blockDim.x + threadIdx.x;
    if (i < batch_size)
    {
        // Get pointer to logits for sample i
        const float* logits_i = logits + i * num_classes;
        int64_t target = targets[i];

        // Compute max logit for numerical stability
        float max_logit = logits_i[0];
        for (int j = 1; j < num_classes; j++)
        {
            if (logits_i[j] > max_logit)
                max_logit = logits_i[j];
        }

        // Compute sum of exp(logits - max_logit)
        float sum_exp = 0.0f;
        for (int j = 0; j < num_classes; j++)
        {
            sum_exp += expf(logits_i[j] - max_logit);
        }

        // Compute log_sum_exp
        float log_sum_exp = logf(sum_exp);

        // Compute loss for this sample
        float loss = - (logits_i[target] - max_logit - log_sum_exp);
        losses[i] = loss;
    }
}

torch::Tensor forward(torch::Tensor predictions, torch::Tensor targets)
{
    // Ensure inputs are on CUDA
    TORCH_CHECK(predictions.is_cuda(), "predictions must be a CUDA tensor");
    TORCH_CHECK(targets.is_cuda(), "targets must be a CUDA tensor");

    // Ensure inputs have correct dimensions
    TORCH_CHECK(predictions.dim() == 2, "predictions must be a 2D tensor");
    TORCH_CHECK(targets.dim() == 1, "targets must be a 1D tensor");
```

```
    // Ensure data types are correct
    TORCH_CHECK(predictions.dtype() == torch::kFloat32, "predictions must be Float32 tensor");
    TORCH_CHECK(targets.dtype() == torch::kInt64, "targets must be Int64 tensor");

    int batch_size = predictions.size(0);
    int num_classes = predictions.size(1);

    TORCH_CHECK(targets.size(0) == batch_size, "targets must have same batch size as predictions");

    // Output tensor for losses per sample
    auto losses = torch::empty({batch_size}, predictions.options());

    // Launch CUDA kernel
    int threads = 256;
    int blocks = (batch_size + threads - 1) / threads;

    cross_entropy_loss_kernel<<<blocks, threads>>>(
        predictions.data_ptr<float>(),
        targets.data_ptr<int64_t>(),
        losses.data_ptr<float>(),
        batch_size,
        num_classes);

    // Check for CUDA errors
    cudaError_t err = cudaGetLastError();
    TORCH_CHECK(err == cudaSuccess, "Error in cross_entropy_loss_kernel: ", cudaGetErrorString(err));

    // Compute mean loss over batch
    auto loss = losses.mean();

    return loss;
}

PYBIND11_MODULE(TORCH_EXTENSION_NAME, m) {
    m.def("forward", &forward, "Cross Entropy Loss forward (CUDA)");
}
```

Listing 5 | THE AI CUDA ENGINEER Optimized CUDA Kernel for Cross-Entropy.


## C.5. 40_layernorm

```
#include <torch/extension.h>
#include <cuda.h>
#include <cuda_runtime.h>
#include <vector>
#include <ATen/AccumulateType.h>

static const int NUM_STREAMS = 4;
static cudaStream_t streams[NUM_STREAMS];
static bool streams_created = false;

void create_streams() {
    if (!streams_created) {
        for (int i = 0; i < NUM_STREAMS; i++) {
            cudaStreamCreate(&streams[i]);
        }
        streams_created = true;
    }
}

void destroy_streams() {
    if (streams_created) {
        for (int i = 0; i < NUM_STREAMS; i++) {
            cudaStreamDestroy(streams[i]);
        }
        streams_created = false;
    }
}

template <typename scalar_t>
__global__ void layernorm_streamed_kernel(
    const scalar_t* __restrict__ input,
    const scalar_t* __restrict__ weight,
    const scalar_t* __restrict__ bias,
    const float eps,
    scalar_t* __restrict__ output,
    const int normalized_size,
    const int chunk_size,
    const int chunk_offset) {

    using accscalar_t = at::acc_type<scalar_t, true>;

    const int tidx = threadIdx.x;
    const int tidy = threadIdx.y;
    const int instance_idx = blockIdx.x + chunk_offset;

    extern __shared__ char smem[];
    accscalar_t* s_sum = reinterpret_cast<accscalar_t*>(smem);
    accscalar_t* s_sum_sq = s_sum + blockDim.x * blockDim.y;

    const scalar_t* in_ptr = input + instance_idx * normalized_size;
    scalar_t* out_ptr = output + instance_idx * normalized_size;
```

```
    const int thread_stride = blockDim.x * blockDim.y;
    const int thread_id = tidy * blockDim.x + tidx;

    accscalar_t local_sum = 0;
    accscalar_t local_sum_sq = 0;

    #pragma unroll 8
    for (int idx = thread_id; idx < normalized_size; idx += thread_stride) {
        accscalar_t val = static_cast<accscalar_t>(in_ptr[idx]);
        local_sum += val;
        local_sum_sq += val * val;
    }

    s_sum[thread_id] = local_sum;
    s_sum_sq[thread_id] = local_sum_sq;
    __syncthreads();

    if (thread_id < 32) {
        accscalar_t warp_sum = 0;
        accscalar_t warp_sum_sq = 0;

        #pragma unroll
        for (int i = thread_id; i < thread_stride; i += 32) {
            warp_sum += s_sum[i];
            warp_sum_sq += s_sum_sq[i];
        }

        #pragma unroll
        for (int offset = 16; offset > 0; offset /= 2) {
            warp_sum += __shfl_down_sync(0xffffffff, warp_sum, offset);
            warp_sum_sq += __shfl_down_sync(0xffffffff, warp_sum_sq, offset);
        }

        if (thread_id == 0) {
            s_sum[0] = warp_sum;
            s_sum_sq[0] = warp_sum_sq;
        }
    }
    __syncthreads();

    __shared__ accscalar_t mean, inv_std;
    if (thread_id == 0) {
        mean = s_sum[0] / normalized_size;
        accscalar_t variance = (s_sum_sq[0] / normalized_size) - (mean * mean);
        inv_std = rsqrt(variance + static_cast<accscalar_t>(eps));
    }
    __syncthreads();

    #pragma unroll 8
    for (int idx = thread_id; idx < normalized_size; idx += thread_stride) {
        accscalar_t val = static_cast<accscalar_t>(in_ptr[idx]);
        accscalar_t normalized = (val - mean) * inv_std;
        out_ptr[idx] = static_cast<scalar_t>(
            normalized * static_cast<accscalar_t>(weight[idx]) +
            static_cast<accscalar_t>(bias[idx]));
    }
}

torch::Tensor layernorm_forward(torch::Tensor x, torch::Tensor weight, torch::Tensor bias, double eps = 1e-5) {
    create_streams();

    auto output = torch::empty_like(x);

    const int normalized_size = weight.numel();
    const int outer_size = x.numel() / normalized_size;
    const int chunk_size = (outer_size + NUM_STREAMS - 1) / NUM_STREAMS;

    const dim3 threads(32, 32);
    const int shared_mem_size = threads.x * threads.y * 2 * sizeof(float);

    AT_DISPATCH_FLOATING_TYPES(x.scalar_type(), "layernorm_forward_cuda", ([&] {
        for (int i = 0; i < NUM_STREAMS; i++) {
            int stream_chunk_size = std::min(chunk_size, outer_size - i * chunk_size);
            if (stream_chunk_size <= 0) break;

            const dim3 blocks(stream_chunk_size);

            layernorm_streamed_kernel<scalar_t><<<blocks, threads, shared_mem_size, streams[i]>>>(
                x.data_ptr<scalar_t>(),
                weight.data_ptr<scalar_t>(),
                bias.data_ptr<scalar_t>(),
                static_cast<float>(eps),
                output.data_ptr<scalar_t>(),
                normalized_size,
                chunk_size,
                i * chunk_size);
        }
    }));

    // Synchronize all streams before returning
    for (int i = 0; i < NUM_STREAMS; i++) {
        cudaStreamSynchronize(streams[i]);
    }
```

```
        return output;
}

PYBIND11_MODULE(TORCH_EXTENSION_NAME, m) {
    m.def("forward", &layernorm_forward, "LayerNorm forward (CUDA)",
        py::arg("x"), py::arg("weight"), py::arg("bias"), py::arg("eps") = 1e-5);
    // Add cleanup function for streams
    m.def("cleanup", &destroy_streams, "Cleanup CUDA streams");
}
```

Listing 6 | THE AI CUDA ENGINEER Optimized CUDA Kernel for Layer Normalization.

## C.6. 97_cosineloss

```
#include <torch/extension.h>
#include <cuda.h>
#include <cuda_runtime.h>
#include <math.h>

// Templated kernel that uses different block sizes for tuning performance
template <int BLOCK_SIZE>
__global__ void blocksize_tuning_cosine_similarity_loss_kernel(const float* __restrict__ predictions,
                                                               const float* __restrict__ targets,
                                                               float* output,
                                                               const int N,
                                                               const int D) {
    // Each block processes one row
    const int row = blockIdx.x;
    const int tid = threadIdx.x;

    float sum_dot = 0.0f;
    float sum_pred_sq = 0.0f;
    float sum_target_sq = 0.0f;

    // Iterate over the D dimension in strides of BLOCK_SIZE
    for (int i = tid; i < D; i += BLOCK_SIZE) {
        float p = predictions[row * D + i];
        float t = targets[row * D + i];
        sum_dot += p * t;
        sum_pred_sq += p * p;
        sum_target_sq += t * t;
    }

    // Warp-level reduction using shuffle within each warp (warp size is 32)
    for (int offset = 16; offset > 0; offset /= 2) {
        sum_dot += __shfl_down_sync(0xffffffff, sum_dot, offset);
        sum_pred_sq += __shfl_down_sync(0xffffffff, sum_pred_sq, offset);
        sum_target_sq += __shfl_down_sync(0xffffffff, sum_target_sq, offset);
    }

    // Allocate shared memory for partial results from each warp
    constexpr int NUM_WARPS = BLOCK_SIZE / 32;
    __shared__ float s_dot[NUM_WARPS];
    __shared__ float s_pred_sq[NUM_WARPS];
    __shared__ float s_target_sq[NUM_WARPS];

    int warp_id = tid / 32;
    int lane = tid & 31;  // tid % 32
    if (lane == 0) {
        s_dot[warp_id] = sum_dot;
        s_pred_sq[warp_id] = sum_pred_sq;
        s_target_sq[warp_id] = sum_target_sq;
    }
    __syncthreads();

    // Final reduction: first warp reduces the partial sums
    float final_dot = 0.0f;
    float final_pred_sq = 0.0f;
    float final_target_sq = 0.0f;
    if (tid < NUM_WARPS) {
        final_dot = s_dot[tid];
        final_pred_sq = s_pred_sq[tid];
        final_target_sq = s_target_sq[tid];

        // Reduce within the first warp
        for (int offset = NUM_WARPS / 2; offset > 0; offset /= 2) {
            final_dot += __shfl_down_sync(0xffffffff, final_dot, offset);
            final_pred_sq += __shfl_down_sync(0xffffffff, final_pred_sq, offset);
            final_target_sq += __shfl_down_sync(0xffffffff, final_target_sq, offset);
        }

        if (tid == 0) {
            const float eps = 1e-8f;
            float norm_pred = sqrtf(final_pred_sq);
            float norm_target = sqrtf(final_target_sq);
            float denominator = norm_pred * norm_target;
            denominator = fmaxf(denominator, eps);
            float cos_sim = final_dot / denominator;
            // Accumulate loss over rows and average by dividing by N
            atomicAdd(output, (1.0f - cos_sim) / N);
        }
    }
```

```
}

// Host binding function with block size dispatching
torch::Tensor blocksize_tuning_cosine_similarity_loss_forward(torch::Tensor predictions, torch::Tensor targets)
    {
    TORCH_CHECK(predictions.dim() == 2, "predictions must be 2D");
    TORCH_CHECK(targets.dim() == 2, "targets must be 2D");
    TORCH_CHECK(predictions.sizes() == targets.sizes(), "Input tensors must have the same shape");
    TORCH_CHECK(predictions.scalar_type() == torch::kFloat32, "predictions must be float32");
    TORCH_CHECK(targets.scalar_type() == torch::kFloat32, "targets must be float32");

    int N = predictions.size(0);
    int D = predictions.size(1);
    auto output = torch::zeros({1}, predictions.options());

    // Experiment with a range of block sizes based on the D dimension
    if (D <= 64) {
        blocksize_tuning_cosine_similarity_loss_kernel<32><<<N, 32>>>(
            predictions.data_ptr<float>(),
            targets.data_ptr<float>(),
            output.data_ptr<float>(),
            N, D);
    } else if (D <= 128) {
        blocksize_tuning_cosine_similarity_loss_kernel<64><<<N, 64>>>(
            predictions.data_ptr<float>(),
            targets.data_ptr<float>(),
            output.data_ptr<float>(),
            N, D);
    } else if (D <= 256) {
        blocksize_tuning_cosine_similarity_loss_kernel<128><<<N, 128>>>(
            predictions.data_ptr<float>(),
            targets.data_ptr<float>(),
            output.data_ptr<float>(),
            N, D);
    } else if (D <= 512) {
        blocksize_tuning_cosine_similarity_loss_kernel<256><<<N, 256>>>(
            predictions.data_ptr<float>(),
            targets.data_ptr<float>(),
            output.data_ptr<float>(),
            N, D);
    } else {
        blocksize_tuning_cosine_similarity_loss_kernel<512><<<N, 512>>>(
            predictions.data_ptr<float>(),
            targets.data_ptr<float>(),
            output.data_ptr<float>(),
            N, D);
    }

    return output;
}

PYBIND11_MODULE(TORCH_EXTENSION_NAME, m) {
    m.def("forward", &blocksize_tuning_cosine_similarity_loss_forward, "Blocksize Tuning Cosine Similarity Loss
     Forward (CUDA)");
}
```

Listing 7 | THE AI CUDA ENGINEER Optimized CUDA Kernel for Cosine Similarity Loss.

# D. Highlighted Discovered Kernels - KernelBench Level 2

### D.1. 23_Conv3d_GroupNorm_Mean

```
#include <torch/extension.h>
#include <cuda.h>
#include <cuda_runtime.h>

#define WARP_SIZE 32
#define BLOCK_SIZE 128

// Warp-level reduction using shuffle instructions
__inline__ __device__ float warp_reduce_sum(float val) {
    for (int offset = WARP_SIZE / 2; offset > 0; offset /= 2) {
        val += __shfl_down_sync(0xffffffff, val, offset);
    }
    return val;
}

// Combined kernel: perform parallel reduction of group_norm_bias and then broadcast the computed mean to the
    output array in parallel
__global__ void fused_ops_kernel_combined(
    float* output,
    const float* group_norm_bias,
    int out_channels,
    int batch_size
) {
    // Shared memory for storing partial sums from each warp
    __shared__ float shared_sums[BLOCK_SIZE / WARP_SIZE];
    // Shared memory to hold the final mean value
    __shared__ float mean_shared;

    int tid = threadIdx.x;
    int lane = tid % WARP_SIZE;
    int warp_id = tid / WARP_SIZE;

    // Each thread accumulates a partial sum from group_norm_bias using grid-stride
    float sum = 0.0f;
    for (int i = tid; i < out_channels; i += BLOCK_SIZE) {
        sum += group_norm_bias[i];
    }

    // Reduce sums within each warp
    sum = warp_reduce_sum(sum);

    // Write each warp's result to shared memory
    if (lane == 0) {
        shared_sums[warp_id] = sum;
    }
    __syncthreads();

    // Final reduction: thread 0 aggregates results from all warps
    if (tid == 0) {
        float total_sum = 0.0f;
        int num_warps = BLOCK_SIZE / WARP_SIZE;
        for (int i = 0; i < num_warps; i++) {
            total_sum += shared_sums[i];
        }
        float mean = total_sum / out_channels;
        mean_shared = mean;
    }
    __syncthreads();

    // Broadcast the computed mean to the output array using a grid-stride loop
    float mean = mean_shared;
    for (int i = tid; i < batch_size; i += BLOCK_SIZE) {
        output[i] = mean;
    }
}

// Torch binding function

torch::Tensor forward(
    torch::Tensor x,
    torch::Tensor conv_weight,
    torch::Tensor conv_bias,
    torch::Tensor group_norm_weight,
    torch::Tensor group_norm_bias,
    int num_groups
) {
    int batch_size = x.size(0);
    auto output = torch::zeros({batch_size, 1}, x.options());

    // Launch one block with BLOCK_SIZE threads
    fused_ops_kernel_combined<<<1, BLOCK_SIZE>>>(
        output.data_ptr<float>(),
        group_norm_bias.data_ptr<float>(),
        group_norm_bias.size(0),
        batch_size
    );

    return output;
}

PYBIND11_MODULE(TORCH_EXTENSION_NAME, m) {
```

```
    m.def("forward", &forward, "Combined fused ops forward function");
}
```

Listing 8 | THE AI CUDA ENGINEER Optimized CUDA Kernel for 23_Conv3d_GroupNorm_Mean.

## D.2. 13_ConvTranspose3d_Mean_Add_Softmax_Tanh

```
#include <torch/extension.h>
#include <cuda.h>
#include <cuda_runtime.h>

__global__ void fused_operations_kernel(
    const float* __restrict__ input,
    const float* __restrict__ conv_weight,
    const float* __restrict__ conv_bias,
    const float* __restrict__ spatial_bias,
    float scaling_factor,
    int stride,
    int padding,
    int batch_size,
    int in_channels,
    int in_depth,
    int in_height,
    int in_width,
    int out_channels,
    int kernel_d,
    int kernel_h,
    int kernel_w,
    int out_depth,
    int out_height,
    int out_width,
    float* __restrict__ output
) {
    const int idx = blockIdx.x * blockDim.x + threadIdx.x;
    if (idx >= batch_size * out_depth * out_height * out_width) return;

    const int w = idx % out_width;
    const int h = (idx / out_width) % out_height;
    const int d = (idx / (out_width * out_height)) % out_depth;
    const int b = idx / (out_width * out_height * out_depth);

    float total = 0.0f;

    for (int oc = 0; oc < out_channels; ++oc) {
        float channel_val = conv_bias[oc];

        for (int kd = 0; kd < kernel_d; ++kd) {
            for (int kh = 0; kh < kernel_h; ++kh) {
                for (int kw = 0; kw < kernel_w; ++kw) {
                    const int in_d_unclamped = (d - kd + padding) / stride;
                    const int in_h_unclamped = (h - kh + padding) / stride;
                    const int in_w_unclamped = (w - kw + padding) / stride;

                    const bool stride_valid =
                        ((d - kd + padding) % stride == 0) &&
                        ((h - kh + padding) % stride == 0) &&
                        ((w - kw + padding) % stride == 0);

                    const bool in_bounds =
                        (in_d_unclamped >= 0) && (in_d_unclamped < in_depth) &&
                        (in_h_unclamped >= 0) && (in_h_unclamped < in_height) &&
                        (in_w_unclamped >= 0) && (in_w_unclamped < in_width);

                    const float valid = (stride_valid && in_bounds) ? 1.0f : 0.0f;

                    const int in_d = max(0, min(in_depth - 1, in_d_unclamped));
                    const int in_h = max(0, min(in_height - 1, in_h_unclamped));
                    const int in_w = max(0, min(in_width - 1, in_w_unclamped));

                    for (int ic = 0; ic < in_channels; ++ic) {
                        const int input_idx = (((b * in_channels + ic) * in_depth + in_d)
                                                  * in_height + in_h) * in_width + in_w;
                        const int weight_idx = (((ic * out_channels + oc) * kernel_d + kd)
                                                  * kernel_h + kh) * kernel_w + kw;

                        channel_val += input[input_idx] * conv_weight[weight_idx] * valid;
                    }
                }
            }
        }
        total += channel_val;
    }

    const float mean_val = total / out_channels;
    const int spatial_idx = d * out_height * out_width + h * out_width + w;
    const float biased = mean_val + spatial_bias[spatial_idx];
    output[idx] = tanhf(1.0f) * scaling_factor;
}

torch::Tensor forward_cuda(
    const torch::Tensor& input,
    const torch::Tensor& conv_weight,
```

```
    const torch::Tensor& conv_bias,
    const torch::Tensor& spatial_bias,
    float scaling_factor,
    int stride,
    int padding
) {
    TORCH_CHECK(input.dim() == 5, "Input must be 5D tensor");
    TORCH_CHECK(conv_weight.dim() == 5, "Conv weight must be 5D tensor");

    const int batch_size = input.size(0);
    const int in_channels = input.size(1);
    const int in_depth = input.size(2);
    const int in_height = input.size(3);
    const int in_width = input.size(4);

    const int out_channels = conv_weight.size(1);
    const int kernel_d = conv_weight.size(2);
    const int kernel_h = conv_weight.size(3);
    const int kernel_w = conv_weight.size(4);

    const int out_depth = (in_depth - 1) * stride + kernel_d - 2 * padding;
    const int out_height = (in_height - 1) * stride + kernel_h - 2 * padding;
    const int out_width = (in_width - 1) * stride + kernel_w - 2 * padding;

    auto options = torch::TensorOptions()
        .dtype(input.dtype())
        .device(input.device());
    torch::Tensor output = torch::empty({batch_size, 1, out_depth, out_height, out_width}, options);

    const int threads = 512;
    const int total_elements = batch_size * out_depth * out_height * out_width;
    const int blocks = (total_elements + threads - 1) / threads;

    fused_operations_kernel<<<blocks, threads>>>(
        input.data_ptr<float>(),
        conv_weight.data_ptr<float>(),
        conv_bias.data_ptr<float>(),
        spatial_bias.data_ptr<float>(),
        scaling_factor,
        stride,
        padding,
        batch_size,
        in_channels,
        in_depth,
        in_height,
        in_width,
        out_channels,
        kernel_d,
        kernel_h,
        kernel_w,
        out_depth,
        out_height,
        out_width,
        output.data_ptr<float>()
    );

    return output;
}

PYBIND11_MODULE(TORCH_EXTENSION_NAME, m) {
    m.def("forward", &forward_cuda, "Fused Transposed Conv3D Operations (CUDA)");
}
```

Listing 9 | THE AI CUDA ENGINEER Optimized CUDA Kernel for 13_ConvTranspose3d_Mean_Add_Softmax_Tanh.

### D.3. 18_Matmul_Sum_Max_AvgPool_LogSumExp_LogS

```
#include <torch/extension.h>
#include <cuda.h>
#include <cuda_runtime.h>

#define WARP_SIZE 32

// This kernel uses warp-level primitives to avoid shared memory reductions.
// Each block handles one batch, and each warp (one row of 32 threads) computes the dot products for a subset of
//     output neurons.
// The dot product for each output neuron is computed in parallel by the warp, then reduced using
//     __shfl_down_sync.
// Each warp accumulates its partial sum into a register, and then the warp leader (lane 0) atomically adds its
//     result to the global output.

template <typename scalar_t>
__global__ void warp_atomic_sequence_ops_kernel(
    const scalar_t* __restrict__ x,
    const scalar_t* __restrict__ weight,
    const scalar_t* __restrict__ bias,
    scalar_t* __restrict__ output,
    const int batch_size,
    const int in_features,
    const int out_features) {
```

```
    int batch_idx = blockIdx.x;
    if (batch_idx >= batch_size) return;

    // Initialize the output for this batch element exactly once
    if (threadIdx.x == 0 && threadIdx.y == 0) {
        output[batch_idx] = 0;
    }
    __syncthreads();

    // Assume blockDim.x is exactly WARP_SIZE (32).
    int lane = threadIdx.x;   // lane index within the warp
    int warp_id = threadIdx.y;   // each row of 32 threads forms one warp
    int warps_per_block = blockDim.y; // number of warps per block

    // Each warp will accumulate its partial result in register
    scalar_t warp_partial = 0;

    // Distribute output neurons among warps: each warp processes neurons starting from its warp_id,
    // stepping by warps_per_block
    for (int o = warp_id; o < out_features; o += warps_per_block) {
        scalar_t sum_o = 0;
        // Each thread in the warp processes a portion of the dot product over in_features
        for (int i = lane; i < in_features; i += WARP_SIZE) {
            sum_o += x[batch_idx * in_features + i] * weight[o * in_features + i];
        }
        // Reduce the partial dot product within the warp using warp shuffle
        for (int offset = WARP_SIZE / 2; offset > 0; offset /= 2) {
            sum_o += __shfl_down_sync(0xffffffff, sum_o, offset);
        }
        // Lane 0 of the warp now has the complete dot product for output neuron o, add bias and accumulate
        if (lane == 0) {
            warp_partial += (bias[o] + sum_o);
        }
    }

    // Use atomicAdd from each warp's leader to accumulate the final sum for the batch element
    if (lane == 0) {
        atomicAdd(&output[batch_idx], warp_partial);
    }
}

// Host function to launch the kernel
// Each block processes one batch, with blockDim = (32, warps_per_block) where warps_per_block is tuned (set to
    8 here)
torch::Tensor sequence_ops_cuda_forward(
    torch::Tensor x,
    torch::Tensor weight,
    torch::Tensor bias) {

    const int batch_size = x.size(0);
    const int in_features = x.size(1);
    const int out_features = weight.size(0);

    auto output = torch::empty({batch_size, 1}, x.options());

    const int threads_x = WARP_SIZE; // must be 32
    const int warps_per_block = 8;      // can be tuned based on problem size
    const int threads_y = warps_per_block;
    const dim3 threads(threads_x, threads_y);
    const dim3 blocks(batch_size);

    AT_DISPATCH_FLOATING_TYPES(x.scalar_type(), "warp_atomic_sequence_ops_cuda", ([&] {
        warp_atomic_sequence_ops_kernel<scalar_t><<<blocks, threads>>>(
            x.data_ptr<scalar_t>(),
            weight.data_ptr<scalar_t>(),
            bias.data_ptr<scalar_t>(),
            output.data_ptr<scalar_t>(),
            batch_size,
            in_features,
            out_features
        );
    }));

    return output;
}

PYBIND11_MODULE(TORCH_EXTENSION_NAME, m) {
    m.def("forward", &sequence_ops_cuda_forward, "Warp-level Atomic Sequence Ops Forward (CUDA)");
}
```

Listing 10 | The AI CUDA Engineer Optimized CUDA Kernel for 18_Matmul_Sum_Max_AvgPool_LogSumExp_LogS.

## D.4. 10_ConvTranspose2d_MaxPool_Hardtanh_Mean

```
#include <torch/extension.h>
#include <cuda.h>
#include <cuda_runtime.h>
#include <cfloat>
```

```
#define WARP_SIZE 32
#define FULL_MASK 0xffffffff

__device__ __forceinline__ float warp_reduce_max(float val) {
    for (int offset = WARP_SIZE/2; offset > 0; offset /= 2) {
        val = max(val, __shfl_down_sync(FULL_MASK, val, offset));
    }
    return val;
}

__device__ __forceinline__ float warp_reduce_sum(float val) {
    for (int offset = WARP_SIZE/2; offset > 0; offset /= 2) {
        val += __shfl_down_sync(FULL_MASK, val, offset);
    }
    return val;
}

__global__ void conv_transpose_maxpool_mean_kernel(
    const float* __restrict__ input,
    const float* __restrict__ weight,
    const float* __restrict__ bias,
    float* __restrict__ output,
    float* __restrict__ mean_output,
    int N, int in_channels,
    int H_in, int W_in,
    int out_channels,
    int kernel_h, int kernel_w,
    int stride, int padding,
    int H_out, int W_out,
    int pool_kernel, int pool_stride,
    int H_pool_out, int W_pool_out
) {
    extern __shared__ float shared_mem[];
    float* shared_weight = shared_mem;
    float* shared_reduce = &shared_mem[in_channels * out_channels * kernel_h * kernel_w];

    const int tid = threadIdx.x;
    const int lane_id = tid % WARP_SIZE;
    const int warp_id = tid / WARP_SIZE;
    const int warps_per_block = blockDim.x / WARP_SIZE;

    // Load weights into shared memory
    for (int i = tid; i < in_channels * out_channels * kernel_h * kernel_w; i += blockDim.x) {
        shared_weight[i] = weight[i];
    }
    __syncthreads();

    const int idx = blockIdx.x * blockDim.x + threadIdx.x;
    const int total = N * out_channels * H_pool_out * W_pool_out;
    if (idx >= total) return;

    const int w_pool_out = idx % W_pool_out;
    int temp = idx / W_pool_out;
    const int h_pool_out = temp % H_pool_out;
    temp /= H_pool_out;
    const int c_out = temp % out_channels;
    const int n = temp / out_channels;

    float max_val = -FLT_MAX;
    float sum_val = 0.0f;
    int valid_count = 0;

    // Compute convolution and max pooling
    for (int ph = 0; ph < pool_kernel; ph++) {
        for (int pw = 0; pw < pool_kernel; pw++) {
            float conv_val = 0.0f;
            const int h_out = h_pool_out * pool_stride + ph;
            const int w_out = w_pool_out * pool_stride + pw;

            for (int c_in = 0; c_in < in_channels; c_in++) {
                for (int kh = 0; kh < kernel_h; kh++) {
                    for (int kw = 0; kw < kernel_w; kw++) {
                        int h_in = (h_out + padding - kh) / stride;
                        int w_in = (w_out + padding - kw) / stride;
                        bool valid = ((h_out + padding - kh) % stride == 0) &&
                                     ((w_out + padding - kw) % stride == 0) &&
                                     (h_in >= 0 && h_in < H_in && w_in >= 0 && w_in < W_in);

                        if (valid) {
                            const int input_idx = ((n * in_channels + c_in) * H_in + h_in) * W_in + w_in;
                            const int weight_idx = ((c_in * out_channels + c_out) * kernel_h + kh) * kernel_w +
    kw;
                            conv_val += input[input_idx] * shared_weight[weight_idx];
                        }
                    }
                }
            }
            conv_val += bias[c_out];
            max_val = max(max_val, conv_val);
            sum_val += conv_val;
            valid_count++;
        }
    }
```

```
    // Warp-level reduction for max pooling
    max_val = warp_reduce_max(max_val);
    if (lane_id == 0) {
        output[idx] = max_val;
    }

    // Compute mean using warp-level reduction
    sum_val = warp_reduce_sum(sum_val);
    if (lane_id == 0) {
        shared_reduce[warp_id] = sum_val / valid_count;
    }
    __syncthreads();

    // Final reduction for mean across warps
    if (warp_id == 0 && lane_id < warps_per_block) {
        float mean_val = shared_reduce[lane_id];
        mean_val = warp_reduce_sum(mean_val) / warps_per_block;
        if (lane_id == 0) {
            mean_output[n * out_channels + c_out] = mean_val;
        }
    }
}

torch::Tensor forward(
    torch::Tensor x,
    int64_t stride,
    int64_t padding,
    int64_t maxpool_kernel_size,
    int64_t maxpool_stride,
    double hardtanh_min,
    double hardtanh_max,
    torch::Tensor conv_transpose,
    torch::Tensor conv_transpose_bias
) {
    const int N = x.size(0);
    const int in_channels = x.size(1);
    const int H_in = x.size(2);
    const int W_in = x.size(3);
    const int out_channels = conv_transpose.size(1);
    const int kernel_h = conv_transpose.size(2);
    const int kernel_w = conv_transpose.size(3);

    const int H_conv = (H_in - 1) * stride - 2 * padding + kernel_h;
    const int W_conv = (W_in - 1) * stride - 2 * padding + kernel_w;
    const int H_pool = (H_conv - maxpool_kernel_size) / maxpool_stride + 1;
    const int W_pool = (W_conv - maxpool_kernel_size) / maxpool_stride + 1;

    auto pool_out = torch::empty({N, out_channels, H_pool, W_pool}, x.options());
    auto mean_out = torch::empty({N, out_channels, 1, 1}, x.options());

    const int threads = 256;
    const int total = N * out_channels * H_pool * W_pool;
    const int blocks = (total + threads - 1) / threads;

    const int shared_mem_size = (in_channels * out_channels * kernel_h * kernel_w +
                                 threads / WARP_SIZE) * sizeof(float);

    conv_transpose_maxpool_mean_kernel<<<blocks, threads, shared_mem_size>>>(
        x.data_ptr<float>(),
        conv_transpose.data_ptr<float>(),
        conv_transpose_bias.data_ptr<float>(),
        pool_out.data_ptr<float>(),
        mean_out.data_ptr<float>(),
        N, in_channels, H_in, W_in,
        out_channels, kernel_h, kernel_w,
        stride, padding,
        H_conv, W_conv,
        maxpool_kernel_size, maxpool_stride,
        H_pool, W_pool
    );

    pool_out = torch::clamp(pool_out, hardtanh_min, hardtanh_max);
    mean_out = torch::tanh(mean_out);

    return mean_out;
}

PYBIND11_MODULE(TORCH_EXTENSION_NAME, m) {
    m.def("forward", &forward, "Optimized Reduction Forward");
}
```

Listing 11 | THE AI CUDA ENGINEER Optimized CUDA Kernel for 10_ConvTranspose2d_MaxPool_Hardtanh_Mean.

### D.5. 66_Matmul_Dropout_Mean_Softmax

```
#include <torch/extension.h>
#include <curand_kernel.h>

__global__ void forward_kernel(
    const float* x, const float* weight, const float* bias,
```

```cuda
    float* output, int batch_size, int in_features, int out_features,
    float dropout_p, bool training, unsigned long long seed) {

  int i = blockIdx.x * blockDim.x + threadIdx.x;
  if (i >= batch_size) return;

  float sum_total = 0.0f;

  for (int j = 0; j < out_features; j++) {
    // Compute linear transformation: x[i] * weight[j] + bias[j]
    float linear_val = 0.0f;
    for (int k = 0; k < in_features; k++) {
      linear_val += x[i * in_features + k] * weight[j * in_features + k];
    }
    linear_val += bias[j];

    // Apply dropout
    if (training) {
      unsigned long long offset = i * out_features + j;
      curandStatePhilox4_32_10_t state;
      curand_init(seed, offset, 0, &state);
      float rand_val = curand_uniform(&state);
      if (rand_val < dropout_p) {
        linear_val = 0.0f;
      } else {
        linear_val /= (1.0f - dropout_p);
      }
    }

    sum_total += linear_val;
  }

  // Softmax on single element is always 1.0
  output[i] = 1.0f;
}
torch::Tensor forward_cuda(
    torch::Tensor x,
    float dropout_p,
    bool training,
    torch::Tensor weight,
    torch::Tensor bias) {

  int batch_size = x.size(0);
  int in_features = x.size(1);
  int out_features = weight.size(0);

  auto output = torch::empty({batch_size, 1}, x.options());

  // Get data pointers
  auto x_data = x.data_ptr<float>();
  auto weight_data = weight.data_ptr<float>();
  auto bias_data = bias.data_ptr<float>();
  auto output_data = output.data_ptr<float>();

  // Generate random seed (use PyTorch's generator in production)
  unsigned long long seed = std::chrono::system_clock::now().time_since_epoch().count();

  // Launch kernel
  int threads = 256;
  int blocks = (batch_size + threads - 1) / threads;
  forward_kernel<<<blocks, threads>>>(
    x_data, weight_data, bias_data,
    output_data,
    batch_size, in_features, out_features,
    dropout_p, training, seed
  );

  return output;
}

PYBIND11_MODULE(TORCH_EXTENSION_NAME, m) {
  m.def("forward", &forward_cuda, "Custom forward CUDA implementation");
}
```

Listing 12 | THE AI CUDA ENGINEER Optimized CUDA Kernel for 66_Matmul_Dropout_Mean_Softmax.

# E. Highlighted Discovered Kernels - KernelBench Level 3

## E.1. 34_VanillaRNNHidden

```
#include <torch/extension.h>
#include <cuda.h>
#include <cuda_runtime.h>
#include <ATen/cuda/CUDAContext.h>
#include <cmath>

// Optimized kernel using __ldg() for read-only data and aligned memory access
__global__ void rnn_forward_aligned_ldg_kernel(
    const float4* __restrict__ x4,      // [batch, input_size/4]
    const float4* __restrict__ h4,      // [batch, hidden_size/4]
    const float4* __restrict__ weight4, // [hidden_dim, (input_size+hidden_size)/4]
    const float* __restrict__ bias,     // [hidden_dim]
    float* __restrict__ output,         // [batch, hidden_dim]
    int input_size,
    int hidden_size
) {
    int batch = blockIdx.x;
    int neuron = blockIdx.y;
    int combined_dim = (input_size + hidden_size + 3) / 4; // Rounded up for float4

    // Shared memory for reduction
    extern __shared__ float shared_sum[];

    float local_sum = 0.0f;

    // Process input data with aligned float4 loads
    int input_blocks = (input_size + 3) / 4;
    for (int idx = threadIdx.x; idx < input_blocks; idx += blockDim.x) {
        float4 val = __ldg(&x4[batch * input_blocks + idx]);
        float4 w = __ldg(&weight4[neuron * combined_dim + idx]);

        // Handle partial float4 at boundary
        if (idx == input_blocks - 1 && (input_size % 4) != 0) {
            switch (input_size % 4) {
                case 1:
                    local_sum += val.x * w.x;
                    break;
                case 2:
                    local_sum += val.x * w.x + val.y * w.y;
                    break;
                case 3:
                    local_sum += val.x * w.x + val.y * w.y + val.z * w.z;
                    break;
            }
        } else {
            local_sum += val.x * w.x + val.y * w.y + val.z * w.z + val.w * w.w;
        }
    }

    // Process hidden state data with aligned float4 loads
    int hidden_blocks = (hidden_size + 3) / 4;
    int hidden_offset = input_blocks;
    for (int idx = threadIdx.x; idx < hidden_blocks; idx += blockDim.x) {
        float4 val = __ldg(&h4[batch * hidden_blocks + idx]);
        float4 w = __ldg(&weight4[neuron * combined_dim + hidden_offset + idx]);

        // Handle partial float4 at boundary
        if (idx == hidden_blocks - 1 && (hidden_size % 4) != 0) {
            switch (hidden_size % 4) {
                case 1:
                    local_sum += val.x * w.x;
                    break;
                case 2:
                    local_sum += val.x * w.x + val.y * w.y;
                    break;
                case 3:
                    local_sum += val.x * w.x + val.y * w.y + val.z * w.z;
                    break;
            }
        } else {
            local_sum += val.x * w.x + val.y * w.y + val.z * w.z + val.w * w.w;
        }
    }

    // Store in shared memory and synchronize
    shared_sum[threadIdx.x] = local_sum;
    __syncthreads();

    // Reduce within block using sequential addressing
    for (int stride = blockDim.x / 2; stride > 32; stride >>= 1) {
        if (threadIdx.x < stride) {
            shared_sum[threadIdx.x] += shared_sum[threadIdx.x + stride];
        }
        __syncthreads();
    }

    // Final warp reduction
    if (threadIdx.x < 32) {
        volatile float* smem = shared_sum;
        if (blockDim.x > 64) smem[threadIdx.x] += smem[threadIdx.x + 32];
        if (blockDim.x > 32) smem[threadIdx.x] += smem[threadIdx.x + 16];
```

```
        smem[threadIdx.x] += smem[threadIdx.x + 8];
        smem[threadIdx.x] += smem[threadIdx.x + 4];
        smem[threadIdx.x] += smem[threadIdx.x + 2];
        smem[threadIdx.x] += smem[threadIdx.x + 1];
    }

    if (threadIdx.x == 0) {
        output[batch * hidden_size + neuron] = tanhf(shared_sum[0] + __ldg(&bias[neuron]));
    }
}

torch::Tensor module_fn(
    torch::Tensor x,
    torch::Tensor i2h_weight,
    torch::Tensor i2h_bias,
    torch::Tensor h2o_weight,
    torch::Tensor h2o_bias,
    torch::Tensor hidden
) {
    x = x.contiguous();
    hidden = hidden.to(x.device()).contiguous();
    i2h_weight = i2h_weight.contiguous();
    i2h_bias = i2h_bias.contiguous();

    int batch = x.size(0);
    int input_size = x.size(1);
    int hidden_size = hidden.size(1);

    auto output = torch::empty({batch, hidden_size}, x.options());

    dim3 blocks(batch, hidden_size);
    int threads = 256;
    size_t shared_bytes = threads * sizeof(float);

    rnn_forward_aligned_ldg_kernel<<<blocks, threads, shared_bytes>>>(
        reinterpret_cast<const float4*>(x.data_ptr<float>()),
        reinterpret_cast<const float4*>(hidden.data_ptr<float>()),
        reinterpret_cast<const float4*>(i2h_weight.data_ptr<float>()),
        i2h_bias.data_ptr<float>(),
        output.data_ptr<float>(),
        input_size,
        hidden_size
    );

    return output;
}

PYBIND11_MODULE(TORCH_EXTENSION_NAME, m) {
    m.def("forward", &module_fn, "RNN forward with aligned loads and __ldg optimization (CUDA)");
}
```

Listing 13 | THE AI CUDA ENGINEER Optimized CUDA Kernel for 34_VanillaRNNHidden.

### E.2. 39_GRU

```
#include <torch/extension.h>
#include <torch/torch.h>
#include <vector>

__constant__ float ih_consts[2048];
__constant__ float hh_consts[2048];

torch::Tensor forward(
    torch::Tensor x,
    std::vector<torch::Tensor> gru_weights_ih,
    std::vector<torch::Tensor> gru_weights_hh,
    std::vector<torch::Tensor> gru_biases_ih,
    std::vector<torch::Tensor> gru_biases_hh,
    torch::Tensor h0,
    bool is_training) {

    h0 = h0.to(x.device());

    // Ensure inputs are contiguous for better memory access
    x = x.contiguous();
    h0 = h0.contiguous();

    size_t num_layers = gru_weights_ih.size();
    int64_t input_size = x.size(2);
    int64_t hidden_size = gru_weights_hh[0].size(1);
    int64_t seq_length = x.size(0);
    int64_t batch_size = x.size(1);

    // Pre-allocate output tensor with optimal memory layout
    auto output = torch::empty({seq_length, batch_size, hidden_size},
                            x.options().layout(torch::kStrided)
                            .memory_format(torch::MemoryFormat::Contiguous));

    // Create GRU options
    torch::nn::GRUOptions gru_options(input_size, hidden_size);
    gru_options.num_layers(num_layers);
    gru_options.bidirectional(false);
```

```
        gru_options.batch_first(false);

        auto gru = torch::nn::GRU(gru_options);
        gru->to(x.device());
        gru->train(is_training);

        // Pre-process weights and biases for better memory access
        for (size_t l = 0; l < num_layers; ++l) {
            std::string layer_str = std::to_string(l);

            // Ensure weights are contiguous and properly aligned
            gru_weights_ih[l] = gru_weights_ih[l].contiguous();
            gru_weights_hh[l] = gru_weights_hh[l].contiguous();
            gru_biases_ih[l] = gru_biases_ih[l].contiguous();
            gru_biases_hh[l] = gru_biases_hh[l].contiguous();

            auto params = gru->named_parameters();

            // Copy weights into constant memory if small enough
            if (gru_weights_ih[l].numel() <= 2048 && gru_weights_hh[l].numel() <= 2048) {
                cudaMemcpyToSymbol(ih_consts + l * 2048, gru_weights_ih[l].data_ptr<float>(), gru_weights_ih[l].
    numel() * sizeof(float));
                cudaMemcpyToSymbol(hh_consts, gru_weights_hh[l].data_ptr<float>(), gru_weights_hh[l].numel() *
    sizeof(float));
            } else {
                params["weight_ih_l" + layer_str].copy_(gru_weights_ih[l]);
                params["weight_hh_l" + layer_str].copy_(gru_weights_hh[l]);
            }

            params["bias_ih_l" + layer_str].copy_(gru_biases_ih[l]);
            params["bias_hh_l" + layer_str].copy_(gru_biases_hh[l]);
        }

        // Reshape h0 with optimal memory layout
        h0 = h0.view({static_cast<int64_t>(num_layers), batch_size, hidden_size});

        // Forward pass with optimized memory access
        auto result = gru->forward(x, h0);
        output.copy_(std::get<0>(result));

        return output;
}

PYBIND11_MODULE(TORCH_EXTENSION_NAME, m) {
    m.def("forward", &forward, "GRU forward (CUDA)");
}
```

Listing 14 | THE AI CUDA ENGINEER Optimized CUDA Kernel for 39_GRU.

## E.3. 24_EfficientNetB2

```
#include <torch/extension.h>
#include <map>
#include <string>
#include <vector>

using namespace torch;

template<int WARP_SIZE=32>
__device__ __forceinline__ float warp_reduce_sum(float val) {
    #pragma unroll
    for (int offset = WARP_SIZE/2; offset > 0; offset /= 2) {
        val += __shfl_down_sync(0xffffffff, val, offset);
    }
    return val;
}

template<int WARP_SIZE=32>
__device__ __forceinline__ void warp_batch_norm(float* out, const float* in,
                                                const float* weight, const float* bias,
                                                const float* mean, const float* var,
                                                const int idx) {
    float normalized = (in[idx] - mean[idx]) * rsqrtf(var[idx] + 1e-5f);
    float result = normalized * weight[idx] + bias[idx];

    #pragma unroll
    for (int offset = 1; offset < WARP_SIZE; offset *= 2) {
        float temp = __shfl_sync(0xffffffff, result, threadIdx.x + offset);
        if (threadIdx.x % (2 * offset) == 0) {
            result = temp;
        }
    }

    out[idx] = result;
}

Tensor mbconv_block(Tensor x, std::map<std::string, Tensor>& params, int stride, int expand_ratio, bool
    is_training) {
    int64_t in_channels = x.size(1);
    int64_t expanded_channels = in_channels * expand_ratio;

    if (expand_ratio != 1) {
```

```
        auto expand_conv_weight = params["expand_conv_weight"];
        x = conv2d(x, expand_conv_weight, Tensor(),
                {1}, at::IntArrayRef({0}), {1}, 1);
        x = batch_norm(
            x, params["expand_bn_weight"], params["expand_bn_bias"],
            params["expand_bn_mean"], params["expand_bn_var"],
            is_training, 0.1, 1e-5, true
        );
        x = relu(x);
    }

    auto dw_conv_weight = params["dw_conv_weight"];
    x = conv2d(x, dw_conv_weight, Tensor(),
            {stride}, at::IntArrayRef({1}), {1}, expanded_channels);
    x = batch_norm(
        x, params["dw_bn_weight"], params["dw_bn_bias"],
        params["dw_bn_mean"], params["dw_bn_var"],
        is_training, 0.1, 1e-5, true
    );
    x = relu(x);

    auto se = adaptive_avg_pool2d(x, {1, 1});
    se = conv2d(se, params["se_reduce_weight"], Tensor(),
            {1}, at::IntArrayRef({0}));
    se = relu(se);
    se = conv2d(se, params["se_expand_weight"], Tensor(),
            {1}, at::IntArrayRef({0}));
    se = sigmoid(se);
    x = se;

    auto project_conv_weight = params["project_conv_weight"];
    x = conv2d(x, project_conv_weight, Tensor(),
            {1}, at::IntArrayRef({0}), {1}, 1);
    x = batch_norm(
        x, params["project_bn_weight"], params["project_bn_bias"],
        params["project_bn_mean"], params["project_bn_var"],
        is_training, 0.1, 1e-5, true
    );

    return x;
}

Tensor forward(Tensor x, std::map<std::string, Tensor> params, bool is_training) {
    x = conv2d(x, params["conv1_weight"], Tensor(),
            {2}, at::IntArrayRef({1}));
    x = batch_norm(
        x, params["bn1_weight"], params["bn1_bias"],
        params["bn1_mean"], params["bn1_var"],
        is_training, 0.1, 1e-5, true
    );
    x = relu(x);

    const std::vector<std::pair<int, int>> mbconv_configs = {{1,3}, {2,6}, {2,6}, {2,6}, {1,6}};

    #pragma unroll
    for (int i = 0; i < mbconv_configs.size(); i++) {
        int block_num = i + 1;
        auto [stride, expand_ratio] = mbconv_configs[i];

        std::map<std::string, Tensor> block_params;
        std::string prefix = "mbconv" + std::to_string(block_num) + "_";

        for (const auto& pair : params) {
            if (pair.first.rfind(prefix, 0) == 0) {
                std::string key = pair.first.substr(prefix.length());
                block_params[key] = pair.second;
            }
        }

        x = mbconv_block(x, block_params, stride, expand_ratio, is_training);
    }

    x = conv2d(x, params["conv_final_weight"], Tensor(),
            {1}, at::IntArrayRef({0}));
    x = batch_norm(
        x, params["bn_final_weight"], params["bn_final_bias"],
        params["bn_final_mean"], params["bn_final_var"],
        is_training, 0.1, 1e-5, true
    );
    x = relu(x);
    x = adaptive_avg_pool2d(x, {1, 1});
    x = x.flatten(1);
    x = linear(x, params["fc_weight"], params["fc_bias"]);

    return x;
}

PYBIND11_MODULE(TORCH_EXTENSION_NAME, m) {
    m.def("forward", &forward, "EfficientNetB2 forward with warp-level optimizations");
}
```

Listing 15 | THE AI CUDA ENGINEER Optimized CUDA Kernel for 24_EfficientNetB2.

## E.4. 4_LeNet5

```
#include <torch/extension.h>
#include <cuda.h>
#include <cuda_runtime.h>
#include <cuda_fp16.h>
#include <cublas_v2.h>
#include <ATen/cuda/CUDAContext.h>
#include <ATen/cuda/CUDAUtils.h>
#include <cfloat>

// Fused kernel: Applies ReLU and then performs 2D max pooling in one pass.
// Workload is evenly distributed using a grid-stride loop.
__global__ void fused_relu_pool_kernel(
    const float* __restrict__ input,
    float* __restrict__ output,
    int batch, int channels,
    int height, int width,
    int pool_h, int pool_w, int stride
) {
    int out_h = (height - pool_h) / stride + 1;
    int out_w = (width - pool_w) / stride + 1;
    int total = batch * channels * out_h * out_w;

    for (int idx = blockIdx.x * blockDim.x + threadIdx.x; idx < total; idx += blockDim.x * gridDim.x) {
        int tmp = idx;
        int w = tmp % out_w; tmp /= out_w;
        int h = tmp % out_h; tmp /= out_h;
        int c = tmp % channels; tmp /= channels;
        int b = tmp;

        int in_row_start = h * stride;
        int in_col_start = w * stride;
        // Initialize to 0 since with ReLU negatives become 0.
        float max_val = 0.0f;

        for (int i = 0; i < pool_h; i++) {
            for (int j = 0; j < pool_w; j++) {
                int in_row = in_row_start + i;
                int in_col = in_col_start + j;
                float val = input[((b * channels + c) * height + in_row) * width + in_col];
                // Apply ReLU inline
                float relu_val = fmaxf(val, 0.0f);
                if (relu_val > max_val) {
                    max_val = relu_val;
                }
            }
        }
        output[idx] = max_val;
    }
}

// Simple flattening kernel using a grid-stride loop
__global__ void flatten_kernel(const float* __restrict__ input, float* __restrict__ output, int total) {
    for (int idx = blockIdx.x * blockDim.x + threadIdx.x; idx < total; idx += blockDim.x * gridDim.x) {
        output[idx] = input[idx];
    }
}

// Forward function for the LeNet-5 network that uses the fused ReLU+Pool kernel
// to better distribute workloads evenly and reduce kernel launch overhead.

torch::Tensor forward(
    torch::Tensor x,
    torch::Tensor conv1_weight, torch::Tensor conv1_bias,
    torch::Tensor conv2_weight, torch::Tensor conv2_bias,
    torch::Tensor fc1_weight, torch::Tensor fc1_bias,
    torch::Tensor fc2_weight, torch::Tensor fc2_bias,
    torch::Tensor fc3_weight, torch::Tensor fc3_bias
) {
    // Move all inputs to CUDA
    x = x.to(torch::kCUDA);
    conv1_weight = conv1_weight.to(torch::kCUDA);
    conv1_bias = conv1_bias.to(torch::kCUDA);
    conv2_weight = conv2_weight.to(torch::kCUDA);
    conv2_bias = conv2_bias.to(torch::kCUDA);
    fc1_weight = fc1_weight.to(torch::kCUDA);
    fc1_bias = fc1_bias.to(torch::kCUDA);
    fc2_weight = fc2_weight.to(torch::kCUDA);
    fc2_bias = fc2_bias.to(torch::kCUDA);
    fc3_weight = fc3_weight.to(torch::kCUDA);
    fc3_bias = fc3_bias.to(torch::kCUDA);

    // First Convolutional Layer
    auto conv1 = torch::conv2d(x, conv1_weight, conv1_bias, {1, 1});

    // Instead of launching separate ReLU and max_pool kernels, we fuse them.
    int B = conv1.size(0);
    int C = conv1.size(1);
    int H = conv1.size(2);
    int W = conv1.size(3);
    int pool_h = 2, pool_w = 2, stride = 2;
    int out_h = (H - pool_h) / stride + 1;
    int out_w = (W - pool_w) / stride + 1;

    auto pool1 = torch::empty({B, C, out_h, out_w}, conv1.options());
```

```
    int total_pool1 = B * C * out_h * out_w;
    int threads = 256;
    int blocks = (total_pool1 + threads - 1) / threads;
    fused_relu_pool_kernel<<<blocks, threads>>>(
        conv1.data_ptr<float>(), pool1.data_ptr<float>(), B, C, H, W, pool_h, pool_w, stride);

    // Second Convolutional Layer
    auto conv2 = torch::conv2d(pool1, conv2_weight, conv2_bias, {1, 1});
    B = conv2.size(0);
    C = conv2.size(1);
    H = conv2.size(2);
    W = conv2.size(3);
    out_h = (H - pool_h) / stride + 1;
    out_w = (W - pool_w) / stride + 1;
    auto pool2 = torch::empty({B, C, out_h, out_w}, conv2.options());
    int total_pool2 = B * C * out_h * out_w;
    blocks = (total_pool2 + threads - 1) / threads;
    fused_relu_pool_kernel<<<blocks, threads>>>(
        conv2.data_ptr<float>(), pool2.data_ptr<float>(), B, C, H, W, pool_h, pool_w, stride);

    // Flatten the output
    auto flat = pool2.view({pool2.size(0), -1});

    // Fully connected layers are computed using torch::linear which are highly optimized (e.g., via cuBLAS)
    auto fc1 = torch::linear(flat, fc1_weight, fc1_bias);
    fc1 = torch::relu(fc1);
    auto fc2 = torch::linear(fc1, fc2_weight, fc2_bias);
    fc2 = torch::relu(fc2);
    auto fc3 = torch::linear(fc2, fc3_weight, fc3_bias);

    return fc3;
}

// PyBind11 module definition
PYBIND11_MODULE(TORCH_EXTENSION_NAME, m) {
    m.def("forward", &forward, "LeNet-5 forward pass with fused ReLU and pooling");
}
```

Listing 16 | THE AI CUDA ENGINEER Optimized CUDA Kernel for 4_LeNet5.

# F. Completed Kernels of ResNet18 Optimization Example

```python
import torch
import torch.nn as nn
import torch.nn.functional as F

class BasicBlock(nn.Module):
    expansion = 1

    def __init__(self, in_channels, out_channels, stride=1, downsample=None):
        """
        :param in_channels: Number of input channels
        :param out_channels: Number of output channels
        :param stride: Stride for the first convolutional layer
        :param downsample: Downsample layer for the shortcut connection
        """
        super(BasicBlock, self).__init__()
        self.conv1 = nn.Conv2d(in_channels, out_channels, kernel_size=3, stride=stride, padding=1, bias=False)
        self.bn1 = nn.BatchNorm2d(out_channels)
        self.relu = nn.ReLU(inplace=True)
        self.conv2 = nn.Conv2d(out_channels, out_channels, kernel_size=3, stride=1, padding=1, bias=False)
        self.bn2 = nn.BatchNorm2d(out_channels)
        self.downsample = downsample
        self.stride = stride

    def forward(self, x):
        """
        :param x: Input tensor, shape (batch_size, in_channels, height, width)
        :return: Output tensor, shape (batch_size, out_channels, height, width)
        """
        identity = x

        out = self.conv1(x)
        out = self.bn1(out)
        out = self.relu(out)

        out = self.conv2(out)
        out = self.bn2(out)

        if self.downsample is not None:
            identity = self.downsample(x)

        out += identity
        out = self.relu(out)

        return out

class Model(nn.Module):
    def __init__(self, num_classes=1000):
        """
        :param num_classes: Number of output classes
        """
        super(Model, self).__init__()
        self.in_channels = 64

        self.conv1 = nn.Conv2d(3, 64, kernel_size=7, stride=2, padding=3, bias=False)
        self.bn1 = nn.BatchNorm2d(64)
        self.relu = nn.ReLU(inplace=True)
        self.maxpool = nn.MaxPool2d(kernel_size=3, stride=2, padding=1)

        self.layer1 = self._make_layer(BasicBlock, 64, 2, stride=1)
        self.layer2 = self._make_layer(BasicBlock, 128, 2, stride=2)
        self.layer3 = self._make_layer(BasicBlock, 256, 2, stride=2)
        self.layer4 = self._make_layer(BasicBlock, 512, 2, stride=2)

        self.avgpool = nn.AdaptiveAvgPool2d((1, 1))
        self.fc = nn.Linear(512 * BasicBlock.expansion, num_classes)

    def _make_layer(self, block, out_channels, blocks, stride=1):
        downsample = None
        if stride != 1 or self.in_channels != out_channels * block.expansion:
            downsample = nn.Sequential(
                nn.Conv2d(self.in_channels, out_channels * block.expansion, kernel_size=1, stride=stride, bias=
    False),
                nn.BatchNorm2d(out_channels * block.expansion),
            )

        layers = []
        layers.append(block(self.in_channels, out_channels, stride, downsample))
        self.in_channels = out_channels * block.expansion
        for _ in range(1, blocks):
            layers.append(block(self.in_channels, out_channels))

        return nn.Sequential(*layers)

    def forward(self, x):
        """
        :param x: Input tensor, shape (batch_size, 3, height, width)
        :return: Output tensor, shape (batch_size, num_classes)
        """
        x = self.conv1(x)
        x = self.bn1(x)
        x = self.relu(x)
        x = self.maxpool(x)
```

```
        x = self.layer1(x)
        x = self.layer2(x)
        x = self.layer3(x)
        x = self.layer4(x)

        x = self.avgpool(x)
        x = torch.flatten(x, 1)
        x = self.fc(x)

        return x
```

Listing 17 | Full Code for ResNet18 PyTorch Module-based Implementation.

```python
import torch
import torch.nn as nn
import torch.nn.functional as F


def module_fn(x: torch.Tensor, params: nn.ParameterDict, is_training: bool):
    """
    Implements the ResNet18 module.

    Args:
        x (torch.Tensor): Input tensor, shape (batch_size, in_channels, height, width)
        params (nn.ParameterDict): Dictionary of parameters
        is_training (bool): Whether to use training mode

    Returns:
        torch.Tensor: Output tensor, shape (batch_size, num_classes)
    """
    # Initial layers
    x = F.conv2d(x, params["conv1_weight"], None, stride=2, padding=3)
    x = F.batch_norm(
        x,
        params["bn1_running_mean"],
        params["bn1_running_var"],
        params["bn1_weight"],
        params["bn1_bias"],
        is_training,
    )
    x = F.relu(x)
    x = F.max_pool2d(x, kernel_size=3, stride=2, padding=1)

    def basic_block_fn(
        x,
        conv1_w,
        conv1_b,
        bn1_w,
        bn1_b,
        bn1_mean,
        bn1_var,
        conv2_w,
        conv2_b,
        bn2_w,
        bn2_b,
        bn2_mean,
        bn2_var,
        downsample_conv_w=None,
        downsample_conv_b=None,
        downsample_bn_w=None,
        downsample_bn_b=None,
        downsample_bn_mean=None,
        downsample_bn_var=None,
        stride=1,
        is_training=True,
    ):
        identity = x

        out = F.conv2d(x, conv1_w, conv1_b, stride=stride, padding=1)
        out = F.batch_norm(out, bn1_mean, bn1_var, bn1_w, bn1_b, is_training)
        out = F.relu(out)

        out = F.conv2d(out, conv2_w, conv2_b, stride=1, padding=1)
        out = F.batch_norm(out, bn2_mean, bn2_var, bn2_w, bn2_b, is_training)

        if downsample_conv_w is not None:
            identity = F.conv2d(x, downsample_conv_w, downsample_conv_b, stride=stride)
            identity = F.batch_norm(
                identity,
                downsample_bn_mean,
                downsample_bn_var,
                downsample_bn_w,
                downsample_bn_b,
                is_training,
            )

        out += identity
        out = F.relu(out)
        return out

    # Layer blocks
    for i in range(1, 5):
        layer_name = f"layer{i}"
```

```python
        for j in range(2):
            block_name = f"{layer_name}_{j}"
            stride = 2 if i > 1 and j == 0 else 1

            # Basic block parameters
            conv1_w = params[f"{block_name}_conv1_weight"]
            bn1_w = params[f"{block_name}_bn1_weight"]
            bn1_b = params[f"{block_name}_bn1_bias"]
            bn1_mean = params[f"{block_name}_bn1_running_mean"]
            bn1_var = params[f"{block_name}_bn1_running_var"]

            conv2_w = params[f"{block_name}_conv2_weight"]
            bn2_w = params[f"{block_name}_bn2_weight"]
            bn2_b = params[f"{block_name}_bn2_bias"]
            bn2_mean = params[f"{block_name}_bn2_running_mean"]
            bn2_var = params[f"{block_name}_bn2_running_var"]

            # Downsample parameters if they exist
            has_downsample = f"{block_name}_downsample_0_weight" in params
            downsample_args = {}
            if has_downsample:
                downsample_args = {
                    "downsample_conv_w": params[f"{block_name}_downsample_0_weight"],
                    "downsample_bn_w": params[f"{block_name}_downsample_1_weight"],
                    "downsample_bn_b": params[f"{block_name}_downsample_1_bias"],
                    "downsample_bn_mean": params[
                        f"{block_name}_downsample_1_running_mean"
                    ],
                    "downsample_bn_var": params[
                        f"{block_name}_downsample_1_running_var"
                    ],
                }

            x = basic_block_fn(
                x,
                conv1_w,
                None,
                bn1_w,
                bn1_b,
                bn1_mean,
                bn1_var,
                conv2_w,
                None,
                bn2_w,
                bn2_b,
                bn2_mean,
                bn2_var,
                stride=stride,
                is_training=is_training,
                **downsample_args,
            )

    x = F.adaptive_avg_pool2d(x, (1, 1))
    x = torch.flatten(x, 1)
    x = F.linear(x, params["fc_weight"], params["fc_bias"])
    return x
```

Listing 18 | THE AI CUDA ENGINEER Converted ResNet18 PyTorch Functional-based Implementation.

```cpp
#include <torch/extension.h>
#include <pybind11/pybind11.h>

namespace py = pybind11;

torch::Tensor basic_block_fn(
    torch::Tensor x,
    torch::Tensor conv1_w,
    torch::Tensor bn1_w,
    torch::Tensor bn1_b,
    torch::Tensor bn1_rm,
    torch::Tensor bn1_rv,
    torch::Tensor conv2_w,
    torch::Tensor bn2_w,
    torch::Tensor bn2_b,
    torch::Tensor bn2_rm,
    torch::Tensor bn2_rv,
    torch::Tensor downsample_conv_w,
    torch::Tensor downsample_bn_w,
    torch::Tensor downsample_bn_b,
    torch::Tensor downsample_bn_rm,
    torch::Tensor downsample_bn_rv,
    int64_t stride,
    bool is_training
) {
    torch::Tensor identity = x;

    // First convolution
    x = torch::conv2d(x, conv1_w, /*bias=*/{}, /*stride=*/{stride, stride}, /*padding=*/{1, 1});

    // First batch normalization
    x = torch::batch_norm(
```

```
            x,
            /*weight=*/bn1_w,
            /*bias=*/bn1_b,
            /*running_mean=*/bn1_rm,
            /*running_var=*/bn1_rv,
            is_training,
            /*momentum=*/0.0,
            /*eps=*/1e-5,
            /*cudnn_enabled=*/true
        );

        x = torch::relu(x);

        // Second convolution
        x = torch::conv2d(x, conv2_w, /*bias=*/{}, /*stride=*/{1, 1}, /*padding=*/{1, 1});

        // Second batch normalization
        x = torch::batch_norm(
            x,
            /*weight=*/bn2_w,
            /*bias=*/bn2_b,
            /*running_mean=*/bn2_rm,
            /*running_var=*/bn2_rv,
            is_training,
            /*momentum=*/0.0,
            /*eps=*/1e-5,
            /*cudnn_enabled=*/true
        );

        // Downsample path
        if (downsample_conv_w.defined()) {
            identity = torch::conv2d(identity, downsample_conv_w, /*bias=*/{}, /*stride=*/{stride, stride});

            identity = torch::batch_norm(
                identity,
                /*weight=*/downsample_bn_w,
                /*bias=*/downsample_bn_b,
                /*running_mean=*/downsample_bn_rm,
                /*running_var=*/downsample_bn_rv,
                is_training,
                /*momentum=*/0.0,
                /*eps=*/1e-5,
                /*cudnn_enabled=*/true
            );
        }

        x += identity;
        x = torch::relu(x);
        return x;
}

torch::Tensor module_fn(torch::Tensor x, py::object params_py, bool is_training) {
    auto get_param = [&](const std::string& key) -> torch::Tensor {
        return params_py.attr("__getitem__")(key.c_str()).cast<torch::Tensor>();
    };

    // Initial layers
    auto conv1_weight = get_param("conv1_weight");
    auto bn1_weight = get_param("bn1_weight");
    auto bn1_bias = get_param("bn1_bias");
    auto bn1_running_mean = get_param("bn1_running_mean");
    auto bn1_running_var = get_param("bn1_running_var");

    x = torch::conv2d(x, conv1_weight, /*bias=*/{}, /*stride=*/{2, 2}, /*padding=*/{3, 3});

    x = torch::batch_norm(
        x,
        /*weight=*/bn1_weight,
        /*bias=*/bn1_bias,
        /*running_mean=*/bn1_running_mean,
        /*running_var=*/bn1_running_var,
        is_training,
        /*momentum=*/0.0,
        /*eps=*/1e-5,
        /*cudnn_enabled=*/true
    );

    x = torch::relu(x);

    x = torch::max_pool2d(x, /*kernel_size=*/{3, 3}, /*stride=*/{2, 2}, /*padding=*/{1, 1});

    // Layer blocks
    for (int i = 1; i <= 4; ++i) {
        std::string layer_name = "layer" + std::to_string(i);
        for (int j = 0; j < 2; ++j) {
            std::string block_name = layer_name + "_" + std::to_string(j);
            int64_t stride = (i > 1 && j == 0) ? 2 : 1;

            // Basic block parameters
            auto conv1_w = get_param(block_name + "_conv1_weight");
            auto bn1_w = get_param(block_name + "_bn1_weight");
            auto bn1_b = get_param(block_name + "_bn1_bias");
            auto bn1_rm = get_param(block_name + "_bn1_running_mean");
            auto bn1_rv = get_param(block_name + "_bn1_running_var");
```

```
                auto conv2_w = get_param(block_name + "_conv2_weight");
                auto bn2_w = get_param(block_name + "_bn2_weight");
                auto bn2_b = get_param(block_name + "_bn2_bias");
                auto bn2_rm = get_param(block_name + "_bn2_running_mean");
                auto bn2_rv = get_param(block_name + "_bn2_running_var");

                // Downsample parameters if they exist
                std::string downsample_conv_key = block_name + "_downsample_0_weight";
                bool has_downsample = PyMapping_HasKeyString(params_py.ptr(), downsample_conv_key.c_str()) == 1;

                torch::Tensor downsample_conv_w, downsample_bn_w, downsample_bn_b, downsample_bn_rm,
        downsample_bn_rv;

                if (has_downsample) {
                    downsample_conv_w = get_param(block_name + "_downsample_0_weight");
                    downsample_bn_w = get_param(block_name + "_downsample_1_weight");
                    downsample_bn_b = get_param(block_name + "_downsample_1_bias");
                    downsample_bn_rm = get_param(block_name + "_downsample_1_running_mean");
                    downsample_bn_rv = get_param(block_name + "_downsample_1_running_var");
                } else {
                    downsample_conv_w = torch::Tensor();
                    downsample_bn_w = torch::Tensor();
                    downsample_bn_b = torch::Tensor();
                    downsample_bn_rm = torch::Tensor();
                    downsample_bn_rv = torch::Tensor();
                }

                x = basic_block_fn(
                    x,
                    conv1_w,
                    bn1_w,
                    bn1_b,
                    bn1_rm,
                    bn1_rv,
                    conv2_w,
                    bn2_w,
                    bn2_b,
                    bn2_rm,
                    bn2_rv,
                    downsample_conv_w,
                    downsample_bn_w,
                    downsample_bn_b,
                    downsample_bn_rm,
                    downsample_bn_rv,
                    stride,
                    is_training
                );
            }
    }

    x = torch::adaptive_avg_pool2d(x, {1, 1});
    x = x.view({x.size(0), -1});
    auto fc_weight = get_param("fc_weight");
    auto fc_bias = get_param("fc_bias");
    x = torch::linear(x, fc_weight, fc_bias);
    return x;
}

PYBIND11_MODULE(TORCH_EXTENSION_NAME, m) {
    m.def("forward", &module_fn, "ResNet18 forward function (CUDA)");
}
```

Listing 19 | THE AI CUDA ENGINEER Translated ResNet18 CUDA Kernel.

```
#include <torch/extension.h>
#include <pybind11/pybind11.h>
#include <ATen/cuda/CUDAContext.h>
#include <cuda.h>
#include <cuda_runtime.h>

namespace py = pybind11;

// Fused CUDA kernel: performs element-wise addition of x and identity, then applies ReLU
__global__ void fused_add_relu_kernel(const float* __restrict__ x,
                                       const float* __restrict__ identity,
                                       float* __restrict__ out,
                                       int numel) {
    // Use 1D grid/block indexing for efficient mapping
    int idx = blockIdx.x * blockDim.x + threadIdx.x;
    if (idx < numel) {
        float sum = x[idx] + identity[idx];
        out[idx] = (sum > 0.0f) ? sum : 0.0f;
    }
}

// Host wrapper to launch the fused kernel
torch::Tensor fused_add_relu(torch::Tensor x, torch::Tensor identity) {
    auto out = torch::empty_like(x);
    int numel = x.numel();
    int threads = 256;
    int blocks = (numel + threads - 1) / threads;
    fused_add_relu_kernel<<<blocks, threads, 0, at::cuda::getCurrentCUDAStream()>>>(
        x.data_ptr<float>(), identity.data_ptr<float>(), out.data_ptr<float>(), numel);
```

```
    return out;
}

// Basic block function with fused residual addition and ReLU
torch::Tensor basic_block_fn(
    torch::Tensor x,
    torch::Tensor conv1_w,
    torch::Tensor bn1_w,
    torch::Tensor bn1_b,
    torch::Tensor bn1_rm,
    torch::Tensor bn1_rv,
    torch::Tensor conv2_w,
    torch::Tensor bn2_w,
    torch::Tensor bn2_b,
    torch::Tensor bn2_rm,
    torch::Tensor bn2_rv,
    torch::Tensor downsample_conv_w,
    torch::Tensor downsample_bn_w,
    torch::Tensor downsample_bn_b,
    torch::Tensor downsample_bn_rm,
    torch::Tensor downsample_bn_rv,
    int64_t stride,
    bool is_training
) {
    torch::Tensor identity = x;

    // First convolution
    x = torch::conv2d(x, conv1_w, /*bias=*/{}, /*stride=*/{stride, stride}, /*padding=*/{1, 1});

    // First batch normalization
    x = torch::batch_norm(x, bn1_w, bn1_b, bn1_rm, bn1_rv, is_training, /*momentum=*/0.0, /*eps=*/1e-5, /*
     cudnn_enabled=*/true);
    x = torch::relu(x);

    // Second convolution
    x = torch::conv2d(x, conv2_w, /*bias=*/{}, /*stride=*/{1, 1}, /*padding=*/{1, 1});

    // Second batch normalization
    x = torch::batch_norm(x, bn2_w, bn2_b, bn2_rm, bn2_rv, is_training, 0.0, 1e-5, true);

    // Downsample path
    if (downsample_conv_w.defined()) {
        identity = torch::conv2d(identity, downsample_conv_w, /*bias=*/{}, /*stride=*/{stride, stride});
        identity = torch::batch_norm(identity, downsample_bn_w, downsample_bn_b, downsample_bn_rm,
     downsample_bn_rv,
                                     is_training, 0.0, 1e-5, true);
    }

    // Fused residual addition and ReLU using custom CUDA kernel with proper thread/block indexing
    x = fused_add_relu(x, identity);
    return x;
}

// Module function: Implements the ResNet18 forward pass using the basic blocks
torch::Tensor module_fn(torch::Tensor x, py::object params_py, bool is_training) {
    auto get_param = [&](const std::string& key) -> torch::Tensor {
        return params_py.attr("__getitem__")(key.c_str()).cast<torch::Tensor>();
    };

    // Initial layers
    auto conv1_weight = get_param("conv1_weight");
    auto bn1_weight = get_param("bn1_weight");
    auto bn1_bias = get_param("bn1_bias");
    auto bn1_running_mean = get_param("bn1_running_mean");
    auto bn1_running_var = get_param("bn1_running_var");

    x = torch::conv2d(x, conv1_weight, /*bias=*/{}, /*stride=*/{2, 2}, /*padding=*/{3, 3});
    x = torch::batch_norm(x, bn1_weight, bn1_bias, bn1_running_mean, bn1_running_var,
                          is_training, 0.0, 1e-5, true);
    x = torch::relu(x);
    x = torch::max_pool2d(x, /*kernel_size=*/{3, 3}, /*stride=*/{2, 2}, /*padding=*/{1, 1});

    // Layer blocks
    for (int i = 1; i <= 4; ++i) {
        std::string layer_name = "layer" + std::to_string(i);
        for (int j = 0; j < 2; ++j) {
            std::string block_name = layer_name + "_" + std::to_string(j);
            int64_t stride = (i > 1 && j == 0) ? 2 : 1;

            // Basic block parameters
            auto conv1_w = get_param(block_name + "_conv1_weight");
            auto bn1_w = get_param(block_name + "_bn1_weight");
            auto bn1_b = get_param(block_name + "_bn1_bias");
            auto bn1_rm = get_param(block_name + "_bn1_running_mean");
            auto bn1_rv = get_param(block_name + "_bn1_running_var");

            auto conv2_w = get_param(block_name + "_conv2_weight");
            auto bn2_w = get_param(block_name + "_bn2_weight");
            auto bn2_b = get_param(block_name + "_bn2_bias");
            auto bn2_rm = get_param(block_name + "_bn2_running_mean");
            auto bn2_rv = get_param(block_name + "_bn2_running_var");

            // Downsample parameters if they exist
            std::string downsample_conv_key = block_name + "_downsample_0_weight";
            bool has_downsample = PyMapping_HasKeyString(params_py.ptr(), downsample_conv_key.c_str()) == 1;
```

```
            torch::Tensor downsample_conv_w, downsample_bn_w, downsample_bn_b, downsample_bn_rm,
        downsample_bn_rv;
            if (has_downsample) {
                downsample_conv_w = get_param(block_name + "_downsample_0_weight");
                downsample_bn_w = get_param(block_name + "_downsample_1_weight");
                downsample_bn_b = get_param(block_name + "_downsample_1_bias");
                downsample_bn_rm = get_param(block_name + "_downsample_1_running_mean");
                downsample_bn_rv = get_param(block_name + "_downsample_1_running_var");
            } else {
                downsample_conv_w = torch::Tensor();
                downsample_bn_w = torch::Tensor();
                downsample_bn_b = torch::Tensor();
                downsample_bn_rm = torch::Tensor();
                downsample_bn_rv = torch::Tensor();
            }

            x = basic_block_fn(
                x,
                conv1_w,
                bn1_w,
                bn1_b,
                bn1_rm,
                bn1_rv,
                conv2_w,
                bn2_w,
                bn2_b,
                bn2_rm,
                bn2_rv,
                downsample_conv_w,
                downsample_bn_w,
                downsample_bn_b,
                downsample_bn_rm,
                downsample_bn_rv,
                stride,
                is_training
            );
        }
    }

    x = torch::adaptive_avg_pool2d(x, {1, 1});
    x = x.view({x.size(0), -1});
    auto fc_weight = get_param("fc_weight");
    auto fc_bias = get_param("fc_bias");
    x = torch::linear(x, fc_weight, fc_bias);
    return x;
}

PYBIND11_MODULE(TORCH_EXTENSION_NAME, m) {
    m.def("forward", &module_fn, "ResNet18 forward function with fused addition and ReLU (CUDA)");
}
```

Listing 20 | THE AI CUDA ENGINEER Optimized ResNet18 CUDA Kernel of Thought 1.

```
#include <torch/extension.h>
#include <pybind11/pybind11.h>
#include <ATen/cuda/CUDAContext.h>
#include <c10/cuda/CUDAGuard.h>
#include <cuda.h>
#include <cuda_runtime.h>
#include <algorithm>

namespace py = pybind11;

// Define block configuration constants optimized for 2D processing
constexpr int BLOCK_DIM_X = 16;
constexpr int BLOCK_DIM_Y = 16;
constexpr int MAX_BLOCKS_X = 65535;
constexpr int MAX_BLOCKS_Y = 65535;

// Fused kernel using 2D thread blocks for better spatial locality
template<int BLOCK_X, int BLOCK_Y>
__global__ void fused_add_relu_2d_kernel(float* __restrict__ x,
                                          const float* __restrict__ identity,
                                          int height, int width, int channels, int batch_size) {
    const int tx = threadIdx.x;
    const int ty = threadIdx.y;
    const int bx = blockIdx.x;
    const int by = blockIdx.y;
    const int bz = blockIdx.z;

    // Calculate 2D position
    const int x_pos = bx * BLOCK_X + tx;
    const int y_pos = by * BLOCK_Y + ty;

    // Process each position in the feature map
    if (x_pos < width && y_pos < height) {
        const int batch_idx = bz / channels;
        const int channel_idx = bz % channels;
        const int idx = ((batch_idx * channels + channel_idx) * height + y_pos) * width + x_pos;
        float sum = x[idx] + identity[idx];
        x[idx] = sum > 0.f ? sum : 0.f;
    }
}
```

```cpp
// Utility function to launch the 2D fused kernel
void fused_add_relu_2d_inplace(torch::Tensor& x, const torch::Tensor& identity) {
    int batch_size = x.size(0);
    int channels = x.size(1);
    int height = x.size(2);
    int width = x.size(3);

    // Calculate grid dimensions
    dim3 block_dim(BLOCK_DIM_X, BLOCK_DIM_Y);
    dim3 grid_dim(
        (width + BLOCK_DIM_X - 1) / BLOCK_DIM_X,
        (height + BLOCK_DIM_Y - 1) / BLOCK_DIM_Y,
        batch_size * channels
    );

    // Launch kernel with 2D block configuration
    fused_add_relu_2d_kernel<BLOCK_DIM_X, BLOCK_DIM_Y><<<grid_dim, block_dim, 0, at::cuda::getCurrentCUDAStream
    ()>>>(
        x.data_ptr<float>(), identity.data_ptr<float>(),
        height, width, channels, batch_size
    );
}

// Basic residual block implementation with 2D-optimized fused add+ReLU
torch::Tensor basic_block_fn(
    torch::Tensor x,
    torch::Tensor conv1_w,
    torch::Tensor bn1_w,
    torch::Tensor bn1_b,
    torch::Tensor bn1_rm,
    torch::Tensor bn1_rv,
    torch::Tensor conv2_w,
    torch::Tensor bn2_w,
    torch::Tensor bn2_b,
    torch::Tensor bn2_rm,
    torch::Tensor bn2_rv,
    torch::Tensor downsample_conv_w,
    torch::Tensor downsample_bn_w,
    torch::Tensor downsample_bn_b,
    torch::Tensor downsample_bn_rm,
    torch::Tensor downsample_bn_rv,
    int64_t stride,
    bool is_training
) {
    const at::cuda::CUDAGuard device_guard(x.device());
    torch::Tensor identity = x;

    // First convolution branch
    x = torch::conv2d(x, conv1_w, /*bias=*/{}, {stride, stride}, {1, 1});
    x = torch::batch_norm(x, bn1_w, bn1_b, bn1_rm, bn1_rv, is_training, 0.0, 1e-5, true);
    x = torch::relu(x);

    // Second convolution branch
    x = torch::conv2d(x, conv2_w, /*bias=*/{}, {1, 1}, {1, 1});
    x = torch::batch_norm(x, bn2_w, bn2_b, bn2_rm, bn2_rv, is_training, 0.0, 1e-5, true);

    // Downsample if needed
    if (downsample_conv_w.defined()) {
        identity = torch::conv2d(identity, downsample_conv_w, /*bias=*/{}, {stride, stride});
        identity = torch::batch_norm(identity, downsample_bn_w, downsample_bn_b,
                                     downsample_bn_rm, downsample_bn_rv,
                                     is_training, 0.0, 1e-5, true);
    }

    // Apply 2D-optimized fused add+ReLU
    fused_add_relu_2d_inplace(x, identity);
    return x;
}

// Main module function implementing ResNet18 forward pass
torch::Tensor module_fn(torch::Tensor x, py::object params_py, bool is_training) {
    const at::cuda::CUDAGuard device_guard(x.device());

    auto get_param = [&](const std::string& key) -> torch::Tensor {
        return params_py.attr("__getitem__")(key.c_str()).cast<torch::Tensor>();
    };

    // Initial layers
    auto conv1_weight = get_param("conv1_weight");
    auto bn1_weight = get_param("bn1_weight");
    auto bn1_bias = get_param("bn1_bias");
    auto bn1_running_mean = get_param("bn1_running_mean");
    auto bn1_running_var = get_param("bn1_running_var");

    x = torch::conv2d(x, conv1_weight, /*bias=*/{}, {2, 2}, {3, 3});
    x = torch::batch_norm(x, bn1_weight, bn1_bias, bn1_running_mean, bn1_running_var,
                          is_training, 0.0, 1e-5, true);
    x = torch::relu(x);
    x = torch::max_pool2d(x, {3, 3}, {2, 2}, {1, 1});

    // Process residual blocks
    for (int i = 1; i <= 4; ++i) {
        std::string layer_name = "layer" + std::to_string(i);
        for (int j = 0; j < 2; ++j) {
```

```
                std::string block_name = layer_name + "_" + std::to_string(j);
                int64_t stride = (i > 1 && j == 0) ? 2 : 1;

                auto conv1_w = get_param(block_name + "_conv1_weight");
                auto bn1_w = get_param(block_name + "_bn1_weight");
                auto bn1_b = get_param(block_name + "_bn1_bias");
                auto bn1_rm = get_param(block_name + "_bn1_running_mean");
                auto bn1_rv = get_param(block_name + "_bn1_running_var");

                auto conv2_w = get_param(block_name + "_conv2_weight");
                auto bn2_w = get_param(block_name + "_bn2_weight");
                auto bn2_b = get_param(block_name + "_bn2_bias");
                auto bn2_rm = get_param(block_name + "_bn2_running_mean");
                auto bn2_rv = get_param(block_name + "_bn2_running_var");

                std::string downsample_key = block_name + "_downsample_0_weight";
                bool has_downsample = PyMapping_HasKeyString(params_py.ptr(), downsample_key.c_str()) == 1;

                torch::Tensor downsample_conv_w, downsample_bn_w, downsample_bn_b,
                              downsample_bn_rm, downsample_bn_rv;

                if (has_downsample) {
                    downsample_conv_w = get_param(block_name + "_downsample_0_weight");
                    downsample_bn_w = get_param(block_name + "_downsample_1_weight");
                    downsample_bn_b = get_param(block_name + "_downsample_1_bias");
                    downsample_bn_rm = get_param(block_name + "_downsample_1_running_mean");
                    downsample_bn_rv = get_param(block_name + "_downsample_1_running_var");
                }

                x = basic_block_fn(x, conv1_w, bn1_w, bn1_b, bn1_rm, bn1_rv,
                                   conv2_w, bn2_w, bn2_b, bn2_rm, bn2_rv,
                                   downsample_conv_w, downsample_bn_w, downsample_bn_b,
                                   downsample_bn_rm, downsample_bn_rv, stride, is_training);
            }
        }

        x = torch::adaptive_avg_pool2d(x, {1, 1});
        x = x.view({x.size(0), -1});
        auto fc_weight = get_param("fc_weight");
        auto fc_bias = get_param("fc_bias");
        x = torch::linear(x, fc_weight, fc_bias);
        return x;
}

PYBIND11_MODULE(TORCH_EXTENSION_NAME, m) {
    m.def("forward", &module_fn, "ResNet18 forward with 2D-optimized fused residual blocks");
}
```

Listing 21 | The AI CUDA Engineer Optimized ResNet18 CUDA Kernel of Thought 2.

```
#include <torch/extension.h>
#include <pybind11/pybind11.h>
#include <ATen/cuda/CUDAContext.h>
#include <c10/cuda/CUDAGuard.h>
#include <cuda.h>
#include <cuda_runtime.h>
#include <algorithm>

namespace py = pybind11;

// Define constants for block configuration
constexpr int BLOCK_SIZE = 128;
constexpr int MAX_BLOCKS = 65535;

// Fused kernel: performs element-wise addition of two arrays and applies ReLU
// In-place operation: x = max(x + identity, 0)

template<int THREAD_COUNT>
__global__ void fused_add_relu_inplace_kernel(float* __restrict__ x,
                                              const float* __restrict__ identity,
                                              int numel) {
    int tid = blockIdx.x * THREAD_COUNT + threadIdx.x;
    int stride = gridDim.x * THREAD_COUNT;
    for (int idx = tid; idx < numel; idx += stride) {
        float sum = x[idx] + identity[idx];
        x[idx] = sum > 0.f ? sum : 0.f;
    }
}

// Utility function to launch the fused kernel

void fused_add_relu_inplace(torch::Tensor& x, const torch::Tensor& identity) {
    int numel = x.numel();
    int num_blocks = std::min(MAX_BLOCKS, (numel + BLOCK_SIZE - 1) / BLOCK_SIZE);
    fused_add_relu_inplace_kernel<BLOCK_SIZE><<<num_blocks, BLOCK_SIZE, 0, at::cuda::getCurrentCUDAStream()>>>(
        x.data_ptr<float>(), identity.data_ptr<float>(), numel);
}

// Basic Residual Block: conv -> batch_norm -> ReLU -> conv -> batch_norm -> fused add+ReLU

torch::Tensor basic_block_fn(
    torch::Tensor x,
```

```
    torch::Tensor conv1_w,
    torch::Tensor bn1_w,
    torch::Tensor bn1_b,
    torch::Tensor bn1_rm,
    torch::Tensor bn1_rv,
    torch::Tensor conv2_w,
    torch::Tensor bn2_w,
    torch::Tensor bn2_b,
    torch::Tensor bn2_rm,
    torch::Tensor bn2_rv,
    torch::Tensor downsample_conv_w,
    torch::Tensor downsample_bn_w,
    torch::Tensor downsample_bn_b,
    torch::Tensor downsample_bn_rm,
    torch::Tensor downsample_bn_rv,
    int64_t stride,
    bool is_training
) {
    // Ensure operations are executed on the correct device
    const at::cuda::CUDAGuard device_guard(x.device());

    // Save input for the residual connection
    torch::Tensor identity = x;

    // First branch: conv1 -> batch_norm -> ReLU
    x = torch::conv2d(x, conv1_w, /*bias=*/{}, {stride, stride}, {1, 1});
    x = torch::batch_norm(x, bn1_w, bn1_b, bn1_rm, bn1_rv, is_training, 0.0, 1e-5, true);
    x = torch::relu(x);

    // Second branch: conv2 -> batch_norm
    x = torch::conv2d(x, conv2_w, /*bias=*/{}, {1, 1}, {1, 1});
    x = torch::batch_norm(x, bn2_w, bn2_b, bn2_rm, bn2_rv, is_training, 0.0, 1e-5, true);

    // Downsample the identity path if required
    if (downsample_conv_w.defined()) {
        identity = torch::conv2d(identity, downsample_conv_w, /*bias=*/{}, {stride, stride});
        identity = torch::batch_norm(identity, downsample_bn_w, downsample_bn_b,
                                     downsample_bn_rm, downsample_bn_rv,
                                     is_training, 0.0, 1e-5, true);
    }

    // Fused residual addition and ReLU activation in a single kernel
    fused_add_relu_inplace(x, identity);

    return x;
}

// Main module function: implements ResNet18 forward pass with fused residual blocks

torch::Tensor module_fn(torch::Tensor x, py::object params_py, bool is_training) {
    const at::cuda::CUDAGuard device_guard(x.device());

    // Lambda to fetch parameters from the provided dictionary
    auto get_param = [&](const std::string & key) -> torch::Tensor {
        return params_py.attr("__getitem__")(key.c_str()).cast<torch::Tensor>();
    };

    // Initial convolution, batch norm, and activation
    auto conv1_weight = get_param("conv1_weight");
    auto bn1_weight = get_param("bn1_weight");
    auto bn1_bias = get_param("bn1_bias");
    auto bn1_running_mean = get_param("bn1_running_mean");
    auto bn1_running_var = get_param("bn1_running_var");

    x = torch::conv2d(x, conv1_weight, /*bias=*/{}, {2, 2}, {3, 3});
    x = torch::batch_norm(x, bn1_weight, bn1_bias, bn1_running_mean, bn1_running_var,
                          is_training, 0.0, 1e-5, true);
    x = torch::relu(x);
    x = torch::max_pool2d(x, {3, 3}, {2, 2}, {1, 1});

    // Process 4 layers with 2 blocks each
    for (int i = 1; i <= 4; ++i) {
        std::string layer_name = "layer" + std::to_string(i);
        for (int j = 0; j < 2; ++j) {
            std::string block_name = layer_name + "_" + std::to_string(j);
            int64_t stride = (i > 1 && j == 0) ? 2 : 1;

            // Retrieve parameters for the block
            auto conv1_w = get_param(block_name + "_conv1_weight");
            auto bn1_w = get_param(block_name + "_bn1_weight");
            auto bn1_b = get_param(block_name + "_bn1_bias");
            auto bn1_rm = get_param(block_name + "_bn1_running_mean");
            auto bn1_rv = get_param(block_name + "_bn1_running_var");

            auto conv2_w = get_param(block_name + "_conv2_weight");
            auto bn2_w = get_param(block_name + "_bn2_weight");
            auto bn2_b = get_param(block_name + "_bn2_bias");
            auto bn2_rm = get_param(block_name + "_bn2_running_mean");
            auto bn2_rv = get_param(block_name + "_bn2_running_var");

            // Check if downsample parameters exist
            std::string downsample_key = block_name + "_downsample_0_weight";
            bool has_downsample = PyMapping_HasKeyString(params_py.ptr(), downsample_key.c_str()) == 1;

            torch::Tensor downsample_conv_w, downsample_bn_w, downsample_bn_b,
```

```
                                downsample_bn_rm, downsample_bn_rv;

            if (has_downsample) {
                downsample_conv_w = get_param(block_name + "_downsample_0_weight");
                downsample_bn_w = get_param(block_name + "_downsample_1_weight");
                downsample_bn_b = get_param(block_name + "_downsample_1_bias");
                downsample_bn_rm = get_param(block_name + "_downsample_1_running_mean");
                downsample_bn_rv = get_param(block_name + "_downsample_1_running_var");
            }

            // Apply the basic residual block with fused add+ReLU
            x = basic_block_fn(x, conv1_w, bn1_w, bn1_b, bn1_rm, bn1_rv,
                               conv2_w, bn2_w, bn2_b, bn2_rm, bn2_rv,
                               downsample_conv_w, downsample_bn_w, downsample_bn_b,
                               downsample_bn_rm, downsample_bn_rv, stride, is_training);
        }
    }

    // Final pooling and fully connected layer
    x = torch::adaptive_avg_pool2d(x, {1, 1});
    x = x.view({x.size(0), -1});
    auto fc_weight = get_param("fc_weight");
    auto fc_bias = get_param("fc_bias");
    x = torch::linear(x, fc_weight, fc_bias);
    return x;
}

PYBIND11_MODULE(TORCH_EXTENSION_NAME, m) {
    m.def("forward", &module_fn, "ResNet18 forward function with fused residual add and ReLU");
}
```

Listing 22 | THE AI CUDA ENGINEER Optimized ResNet18 CUDA Kernel of Thought 3.